

542931

(12)特許協力条約に基づいて公開された国際出願

(19) 世界知的所有権機関
国際事務局



(43) 国際公開日
2004 年12 月23 日 (23.12.2004)

PCT

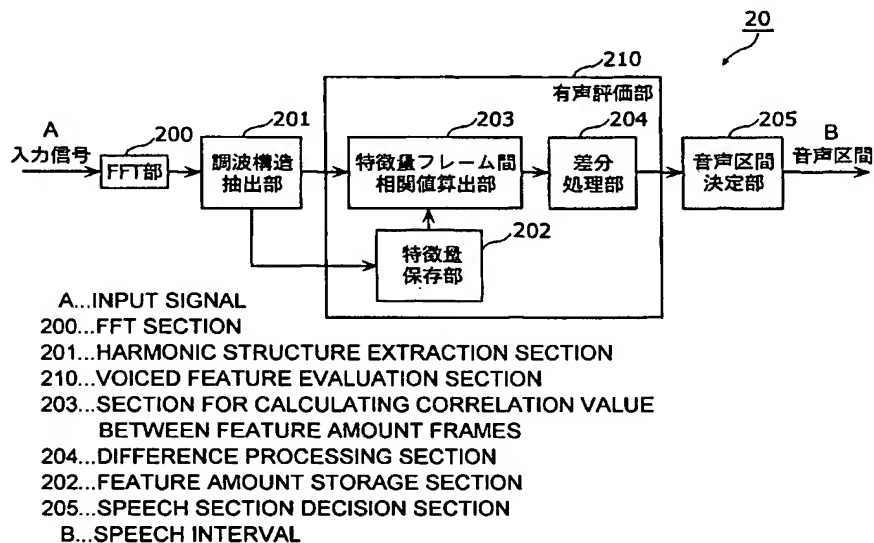
(10) 国際公開番号
WO 2004/111996 A1

- (51) 国際特許分類⁷: G10L 11/02, 11/06 (72) 発明者; および
(21) 国際出願番号: PCT/JP2004/008051 (75) 発明者/出願人 (米国についてのみ): 鈴木 哲 (SUZUKI, Tetsu). 金森 丈郎 (KANAMORI, Takeo). 河村 岳 (KAWAMURA, Takashi).
(22) 国際出願日: 2004 年6 月3 日 (03.06.2004) (74) 代理人: 新居 広守 (NII, Hiromori); 〒5320011 大阪府大阪市淀川区西中島3丁目11番26号 新大阪末広センタービル3F 新居国際特許事務所内 Osaka (JP).
(25) 国際出願の言語: 日本語 (81) 指定国 (表示のない限り、全ての種類の国内保護が可能): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BW, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NA, NL, NO, NZ, OM, PG, PH, PL, PT, RO, RU, SC, SD, SE,
(26) 国際公開の言語: 日本語
(30) 優先権データ: 特願2003-165946 2003 年6 月11 日 (11.06.2003) JP
(71) 出願人 (米国を除く全ての指定国について): 松下電器産業株式会社 (MATSUSHITA ELECTRIC INDUSTRIAL CO., LTD.) [JP/JP]; 〒5718501 大阪府門真市大字門真1006番地 Osaka (JP).

[続葉有]

(54) Title: ACOUSTIC INTERVAL DETECTION METHOD AND DEVICE

(54) 発明の名称: 音響区間検出方法および装置



(57) Abstract: There is provided a harmonic-structured acoustic signal interval detection device not depending on the level fluctuation of the input signal, having an excellent real time characteristic, and noise resistance. The device includes: an FFT section (200) for subjecting the input signal to FFT and calculating a power spectrum component for each frame; a harmonic structure extraction section (201) for leaving only the harmonic structure from the power spectrum component; a voiced feature evaluation section (210) for evaluating correlation between the frames of harmonic structure extracted in the harmonic structure extraction section (201), thereby evaluating whether the interval is a vowel interval and extracting the voiced feature interval; and a speech interval decision section (205) for deciding a speech interval according to the continuity and durability of the output of the voiced feature evaluation section (210).

(57) 要約: 入力信号のレベル変動に依存せず、リアルタイム性に優れ、耐雑音性のある調波構造性音響信号区間検出装置は、入力信号にFFTを施し、フレームごとにパワースペクトル成分を求めるFFT部(200)と、パワースペクトル成分より調波構造のみを残す調波構造抽出部(201)

[続葉有]

WO 2004/111996 A1



SG, SK, SL, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US,
UZ, VC, VN, YU, ZA, ZM, ZW.

BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN,
TD, TG).

- (84) 指定国 (表示のない限り、全ての種類の広域保護が可能): ARIPO (BW, GH, GM, KE, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG, ZM, ZW), ユーラシア (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), ヨーロッパ (AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IT, LU, MC, NL, PL, PT, RO, SE, SI, SK, TR), OAPI (BF,

添付公開書類:

- 国際調査報告書
- 補正書・説明書

2文字コード及び他の略語については、定期発行される各PCTガゼットの巻頭に掲載されている「コードと略語のガイダンスノート」を参照。

明 細 書

音響区間検出方法および装置

5 技術分野

本発明は、入力音響信号から調波構造を有する信号とくに音声が含まれる区間を音声区間として検出する調波構造信号区間および調波構造型音響信号区間検出方法に関し、特に、環境雑音下における調波構造信号および調波構造型音響信号区間検出方法に関する。

10

背景技術

人間の音声は、声帯の振動と発声器官の共振によって形成されており、音の大きさや音の高低を区別するために声帯を制御して振動の周波数を変化させたり、鼻や舌などの発声器官の位置つまり声道形状を変動させたりすることで、人はさまざまな音を発声していることが知られている。このように生成される音声を、音響信号として捕えると、その特徴は、周波数とともに緩やかに変化する成分である、スペクトル包絡と、短時間の周期的（有声母音などの場合）にまたは非周期的に変化する成分（子音や無声母音の場合）である、スペクトル微細構造から構成されていることが知られている。前者のスペクトル包絡成分が発声器官の共振特性を表しており、人間の喉や口の形をあらわす特徴量として用いられ、たとえば音声認識の特徴量としても用いられている。一方、後者のスペクトル微細構造は、音源の周期性を表しており、声帯の基本周期（ピッチ）、音の高低を表す特徴量として用いられている。音声信号のスペクトルは、これら2つの要素の積で表現されている。とくに母音部などにおいて、後者の基本周期およびその高調波成分をよく残している信号は、音声の

調波構造とも呼ばれている。

従来、入力音響信号から音声区間を検出する手法は、様々提案されている。それらを大きく分類すると、入力音響信号の帯域パワーやスペクトルの概形を示すスペクトル包絡などの振幅情報を用いて識別する方法
5 (以下、「方法 1」という。)、口映像を動画像解析することにより、その開閉を検出する方法(以下、「方法 2」という。)、音声や雑音を表現する音響モデルと入力音響信号の音響特徴量とを比較することにより音声区間を検出する方法(以下、「方法 3」という。)、および音声の調音器官の特徴である声道形状によって形成されるスペクトル包絡形状や声帯振動
10 によって形成される調波構造に着目して音声区間を決定する方法(以下、「方法 4」という。) などがある。

しかし、方法 1 では、もともと振幅情報だけで音声と雑音とを識別することが難しいという問題を含んでいる。このため、方法 1 では、音声区間と雑音区間とを仮定し、音声区間と雑音区間とを区別するために設
15 定したしきい値を再学習することにより、音声区間の検出を行なっている。したがって、学習過程において雑音区間の振幅が音声区間の振幅に対して大きくなる(すなわち音声雑音比(以下、「SNR」という。))が 0 dB 程度まで低下する) と、雑音区間であるか音声区間であるかの仮定そのものの精度が性能に影響し、しきい値学習の精度が劣化してしまう。
20 う。その結果として、音声区間検出の性能が劣化するという問題がある。

また、方法 2 では、例えば音入力を用いずに画像だけを用いて口が開いたことを検出するようにすれば、その音声区間検出推定精度は、SNR とは無関係に一定に保つことが可能である。しかし、画像解析処理は音声信号の解析処理に比べて、コストが高いことと、口がカメラの方向
25 に向いていない場合には音声区間の検出ができないという問題がある。

さらに、方法 3 では、想定した環境雑音下での性能は確保されるもの

の、雑音を想定することそのものが難しいため、この方法を使用できる環境は限定的となってしまう。その場の雑音環境を学習する手法も提案されているが、振幅情報を利用する方法（方法１）と同様に、学習方法の精度に依存して性能が劣化するという問題もある。

- 5 一方、音声の調音器官の特徴である、声道形状によって形成されるスペクトル包絡形状や声帯振動によって形成される調波構造に着目して音声区間を決定する方法（方法４）も提案されてきた。

スペクトル包絡形状を利用した方法には、帯域パワー例えばケプストラムの連続性を評価する方法などがあるが、SNRが低下した状況では
10 雑音のオフセット成分との区別がつきにくくなるため、性能が劣化する。

調波構造に着目した方法として、ピッチ検出法はその手法の一つであり、時間軸上の自己相関や高次ケフレンシーを抽出する方法、周波数軸上の自己相関を行なう方法等が提案されている。しかし、これらの方法は、対象とする信号が単一のピッチ（高調波の基本周波数）を持つ信号
15 でない場合には音声区間の抽出が困難であり、環境雑音によって抽出誤りが発生し易い等の問題がある。

また、複数種類の音響信号が混在した音響信号から、人の音声や特定の楽器音等の調波構造を持った音響信号を強調したり、抑圧したり、分離抽出したりする技術が知られている。例えば音声信号に対しては、雑
20 音と音声信号とが混在した音響信号から雑音のみを抑圧する雑音抑圧装置（たとえば、特開平９－１５３７６９号公報参照。）が、また音楽に対しては演奏に含まれる旋律の分離方法や除去方法（たとえば、特開平１
1－１４３４６０号公報参照。）が、それぞれ提案されている。

しかし、特開平９－１５３７６９号公報に記載の方法では、入力信号
25 の線形予測残差信号を帯域ごとに観察することで音声および非音声の検出を行っている。したがって、線形予測がうまく機能しない低SNRの

非定常雑音下では性能が劣化するという問題がある。

また、特開平 1 1 - 1 4 3 4 6 0 号公報に記載の方法は、同一の音程の音が一定時間持続するという音楽の旋律特有の性質を利用した方法である。このため、この方法を、音声と雑音との区別にそのまま用いることは困難であるという問題がある。音響の分離や除去を目的としない場合には、その処理量の多さが問題となる。

調波構造を表現する音響特徴量そのものを評価関数に用いる手法（たとえば、特開 2 0 0 1 - 2 2 2 2 8 9 号公報参照。）も提案されている。図 3 2 は、特開 2 0 0 1 - 2 2 2 2 8 9 号公報で提案されている方法を用いた音声区間決定装置の概略構成を示すブロック図である。

図 3 2 に示される音声区間検出装置 1 0 は、入力信号中の音声区間を決定する装置であり、F F T（Fast Fourier Transform）部 1 0 0 と、調波構造評価部 1 0 1 と、調波構造ピーク検出部 1 0 2 と、ピッチ候補検出部 1 0 3 と、フレーム間振幅差分調波構造評価部 1 0 4 と、音声区間決定部 1 0 5 とを備える。

F F T 部 1 0 0 は、入力信号に対し、フレーム（たとえば、1 フレームは、1 0 m s e c）ごとに F F T 処理を行ない、入力信号を周波数変換し、各種の分析を行なう。調波構造評価部 1 0 1 は、F F T 部 1 0 0 より得られた周波数分析結果より、フレームごとに調波構造を有するかどうかの評価を行なう。調波構造ピーク検出部 1 0 2 は、調波構造評価部 1 0 1 で抽出された調波構造をローカルピーク形状に変換し、ローカルピークを検出する。

ピッチ候補検出部 1 0 3 は、調波構造ピーク検出部 1 0 2 で検出されたローカルピークを時間軸方向（フレーム方向）にトラッキングすることによりピッチ検出を行なう。ピッチとは、調波構造の基本周波数のことである。

フレーム間振幅差分調波構造評価部 104 は、FFT 部 100 における周波数分析の結果得られた振幅をフレーム間で差分し、差分値を求め、その差分値より着目しているフレームが調波構造を有するか否かの評価を行なう。

- 5 音声区間決定部 105 は、ピッチ候補検出部 103 で検出されたピッチと、フレーム間振幅差分調波構造評価部 104 の評価結果とを総合的に判断し、音声区間を決定する。

したがって、図 32 に示される音声区間検出装置 10 では、単一のピッチのみを有する音響信号のみならず、複数のピッチを有する音響信号
10 であっても、音声区間を決定できる。

しかしながら、ピッチ候補検出部 103 において、ローカルピークをトラッキングする際には、ローカルピークの出現や消滅などを考慮しなければならず、これらを考慮しつつ、高精度でピッチを検出するのは困難である。

- 15 また、ピークという極大値を扱う性質上、雑音に対する耐性もあまり期待できない。さらに、時間的な変動を評価するために、フレーム間振幅差分調波構造評価部 104 においては、フレーム間差分に対して調波構造の有無を評価しているが、単に、振幅の差分を用いているため、調波構造の有する情報が失われてしまうだけではなく、例えば突発雑音が
20 生じた場合には、差分値として突発雑音の音響特徴量がそのまま評価されてしまうという問題がある。

そこで、本発明は上述の課題を解決するためになされたものであり、入力信号のレベル変動に依存せず、精度良く音声区間を検出可能な調波構造型音響信号区間検出方法および装置を提供することを目的とする。

- 25 また、リアルタイム性に優れた調波構造型音響信号区間検出方法および装置を提供することも目的とする。

発明の開示

本発明のある局面に係る調波構造型音響信号区間検出方法は、入力音響信号から調波構造を有する信号とくに音声が含まれる区間を音声区間として検出する調波構造型音響信号区間検出方法であって、前記入力音響信号に対し、所定の時間で区切られたフレーム単位で音響特徴量を抽出する音響特徴量抽出ステップと、前記音響特徴量の持続性を評価し、評価結果に従って音声区間を決定する区間決定ステップとを含むことを特徴とする。

このように、音響特徴量の持続性を評価することにより、音声区間の決定を行なっている。このため、ローカルピークをトラッキングする従来の方法のようにローカルピークの出現や消滅など、入力信号のレベル変動を考慮する必要がなく、精度よく音声区間を決定することができる。

好ましくは、前記音響特徴量抽出ステップでは、前記入力音響信号に対しフレーム単位で周波数変換を行ない、前記周波数変換の結果より調波構造のみを強調し、前記音響特徴量を抽出することを特徴とする。

音声（特に母音）には、調波構造が見られる。このため、調波構造を強調した音響特徴量を用いて音声区間を決定することにより、さらに精度よく音声区間を決定することができる。

さらに好ましくは、前記音響特徴量抽出ステップでは、さらに、前記周波数変換の結果より調波構造を抽出し、当該調波構造を含む所定の帯域の周波数変換の結果を、前記音響特徴量とすることを特徴とする。

調波構造が保たれている帯域のみからなる音響特徴量を用いて音声区間を決定することにより、さらに精度よく音声区間を決定することができる。

さらに好ましくは、前記区間決定ステップでは、前記音響特徴量のフレーム間における相関値に基づいて、前記持続性を評価することを特徴

とする。

このように、調波構造の持続性をフレーム間の音響特徴量の相関値により評価している。このため、フレーム間での振幅差分を取り調波構造の持続性を評価する従来方法に比べ、調波構造の有する情報を残した評価が可能である。よって、短いフレームにわたる突発雑音が生じたような場合であっても、そのような突発雑音を音声区間として検出することがなくなり、精度よく音声区間を決定することができる。

さらに好ましくは、前記区間決定ステップは、前記音響特徴量の持続性を評価する評価値を算出する評価ステップと、前記評価値の時間的な連続性を評価し、評価結果に従って音声区間を決定する音声区間決定ステップとを含むことを特徴とする。

音声区間決定ステップでの処理は、実施の形態に述べるように、時間的に連続する有声区間（評価値のみから求められた音声区間）を連結して音声区間を検出する処理に相当する。このように、時間的に連続する有声区間を連結し、音声区間を決定することにより、母音に比べ調波構造性評価値が小さい子音をも音声区間と決定することができる。

さらに、調波構造を有する区間を、詳細に評価することにより、音声か非音声である音楽かどうかを判定することが可能である。調波構造を有すると判定されたフレームにおいて、フレーム内部で最大あるいは最小の調波構造性値が検出された帯域の番号指数を連続的に評価することで、その検出が可能である。

また、フレーム間における調波構造持続性評価値を用いて、調波構造があるとみなされた区間において、該評価値の分散を用いて、音声あるいは音楽など調波構造が持続した区間からの変移なのか、調波構造を持つ突発的なノイズなのかを判別することが可能である。

また、上記調波構造に関する特徴を有する区間以外の区間に対しては、

無音とみなせるほど入力信号が小さい区間あるいは調波構造を有しない非調波構造の区間を判定することができる。

また、実施の形態 5 で示すように、音入力しながらフレーム単位で調波構造性の判定を行なう方法を開示する。

- 5 さらに好ましくは、前記区間決定ステップは、さらに、所定数のフレームにわたる前記評価ステップにおいて算出される前記評価値と第 1 の所定しきい値との比較に基づいて、前記入力音響信号の音声雑音比を推定するステップと、推定された前記音声雑音比が第 2 の所定しきい値以上の場合には、前記評価ステップにおいて算出される前記評価値に基づいて前記音声区間を決定するステップとを含み、前記音声区間決定ステップでは、前記音声雑音比が前記第 2 の所定しきい値未満の場合に、前記評価値の時間的な連続性を評価し、評価結果に従って前記音声区間を決定することを特徴とする。
- 10

- これにより、入力音響信号の推定音声雑音比が良好な場合には、音響
15 特徴量の持続性を評価する評価値の時間的な連続性を評価し、前記音声区間を決定する処理を省略することができる。このため、リアルタイム性に優れた音声区間の検出が可能になる。

- なお、本発明は、以上のような調波構造性音響信号区間検出方法として実現することができるだけでなく、そのステップを手段とする調波構造性音響信号区間検出装置として実現したり、調波構造性音響信号区間
20 検出方法の各ステップをコンピュータに実行させるためのプログラムとして実現したりすることもできる。そのようなプログラムは、CD-ROM等の記録媒体やインターネット等の伝送媒体を介して配信することができるのはいうまでもない。

- 25 以上のように、本発明に係る調波構造性音響信号区間検出方法および装置によると、音声区間と雑音区間との精度良い選別が可能となり、特

に、音声認識方法の前処理として本発明を適用することにより、音声認識率を向上させることができ、その実用的価値は極めて高い。また、IC (Integrated Circuit) レコーダなどに使用することにより音声区間のみを録音したりすることにより、記録容量の効率利用も可能である。

5

図面の簡単な説明

図 1 は、本発明の実施の形態 1 に係る音声区間検出装置のハードウェア構成を示すブロック図である。

図 2 は、実施の形態 1 に係る音声区間検出装置が実行する処理のフローチャートである。

10

図 3 は、調波構造抽出部による調波構造抽出処理のフローチャートである。

図 4 (a) ~ 図 4 (f) は、各フレームにおけるスペクトル成分から調波構造のみを残したスペクトル成分を抽出する過程を模式的に示す図である。

15

図 5 (a) ~ 図 5 (f) は、本発明による入力信号の変換の遷移を示す図である。

図 6 は、音声区間決定処理のフローチャートである。

図 7 は、本発明の実施の形態 2 に係る音声区間検出装置のハードウェア構成を示すブロック図である。

20

図 8 は、実施の形態 2 に係る音声区間検出装置が実行する処理のフローチャートである。

図 9 は、実施の形態 3 に係る音声区間検出装置のハードウェア構成を示すブロック図である。

図 10 は、音声区間検出装置が実行する処理のフローチャートである。

25

図 11 は、調波構造抽出処理を説明するための図である。

図 1 2 は、調波構造抽出処理の詳細を示すフローチャートである。

図 1 3 (a) は、入力信号のパワースペクトルを示す図である。図 1 3 (b) は、調波構造型値 $R(i)$ を示す図である。図 1 3 (c) は帯域番号 $N(i)$ を示す図である。図 1 3 (d) は重み付き帯域番号 $N_e(i)$ を示す図である。図 1 3 (e) は補正調波構造型値 $R'(i)$ を示す図である。

図 1 4 (a) は、入力信号のパワースペクトルを示す図である。図 1 4 (b) は、調波構造型値 $R(i)$ を示す図である。図 1 4 (c) は帯域番号 $N(i)$ を示す図である。図 1 4 (d) は重み付き帯域番号 $N_e(i)$ を示す図である。図 1 4 (e) は補正調波構造型値 $R'(i)$ を示す図である。

図 1 5 (a) は、入力信号のパワースペクトルを示す図である。図 1 5 (b) は、調波構造型値 $R(i)$ を示す図である。図 1 5 (c) は帯域番号 $N(i)$ を示す図である。図 1 5 (d) は重み付き帯域番号 $N_e(i)$ を示す図である。図 1 5 (e) は補正調波構造型値 $R'(i)$ を示す図である。

図 1 6 (a) は、入力信号のパワースペクトルを示す図である。図 1 6 (b) は、調波構造型値 $R(i)$ を示す図である。図 1 6 (c) は帯域番号 $N(i)$ を示す図である。図 1 6 (d) は重み付き帯域番号 $N_e(i)$ を示す図である。図 1 6 (e) は補正調波構造型値 $R'(i)$ を示す図である。

図 1 7 は、音声・音楽区間決定処理の詳細なフローチャートである。

図 1 8 は、実施の形態 4 に係る音声区間検出装置のハードウェア構成を示すブロック図である。

図 1 9 は、音声区間検出装置が実行する処理のフローチャートである。

図 2 0 は、調波構造抽出処理の詳細を示すフローチャートである。

図 2 1 は、音声区間決定処理の詳細を示すフローチャートである。

図 2 2 (a) は入力信号のパワースペクトルを示す図である。図 2 2 (b) は調波構造成値 $R(i)$ を示す図である。図 2 2 (c) は、重み付き分散 $V_e(i)$ を示す図である。図 2 2 (d) は連結前の音声区間を示す図である。図 2 2 (e) は連結後の音声区間を示す図である。

図 2 3 (a) は入力信号のパワースペクトルを示す図である。図 2 3 (b) は調波構造成値 $R(i)$ を示す図である。図 2 3 (c) は、重み付き分散 $V_e(i)$ を示す図である。図 2 3 (d) は連結前の音声区間を示す図である。図 2 3 (e) は連結後の音声区間を示す図である。

10 図 2 4 は、調波構造抽出処理の他の一例を示すフローチャートである。

図 2 5 (a) は入力信号を示す図である。図 2 5 (b) は入力信号のパワースペクトルを示す図である。図 2 5 (c) は調波構造成値 $R(i)$ を示す図である。図 2 5 (d) は重み付き調波構造成値 $R_e(i)$ を示す図である。図 2 5 (e) は補正調波構造成値 $R'(i)$ を示す図である。

15 図 2 6 (a) は入力信号を示す図である。図 2 6 (b) は入力信号のパワースペクトルを示す図である。図 2 6 (c) は調波構造成値 $R(i)$ を示す図である。図 2 6 (d) は重み付き調波構造成値 $R_e(i)$ を示す図である。図 2 6 (e) は補正調波構造成値 $R'(i)$ を示す図である。

図 2 7 は、実施の形態 5 に係る音声区間検出装置 60 の構成を示すブロック図である。

図 2 8 は、音声区間検出装置の実行する処理のフローチャートである。

図 2 9 (a) ~ 図 2 9 (d) は、調波構造成区間の連結を説明するための図である。

図 3 0 は、調波構造成フレーム仮判定処理の詳細なフローチャートである。

図 3 1 は、調波構造成区間確定処理の詳細なフローチャートである。

図 3 2 は、従来の音声区間決定装置の概略のハードウェア構成を示す図である。

発明を実施するための最良の形態

5 (実施の形態 1)

以下、図面を参照しながら本発明の実施の形態 1 に係る音声区間検出装置について説明する。図 1 は、本実施の形態に係る音声区間検出装置 2 0 のハードウェア構成を示すブロック図である。

音声区間検出装置 2 0 は、入力音響信号（以下、単に「入力信号」という。）の中から人間が発声している区間である音声区間を決定する装置であり、FFT 部 2 0 0 と、調波構造抽出部 2 0 1 と、有声評価部 2 1 0 と、音声区間決定部 2 0 5 とを備える。

FFT 部 2 0 0 は、入力信号に FFT を施し、フレームごとにパワースペクトル成分を求める。ここで、1 フレームあたりの時間は 1 0 m s e c とするが、この時間に限定されるものではない。

調波構造抽出部 2 0 1 は、FFT 部 2 0 0 で抽出されたパワースペクトル成分から雑音成分等を取り除き、調波構造のみを残したパワースペクトル成分を抽出する。

有声評価部 2 1 0 は、調波構造抽出部 2 0 1 で抽出された調波構造のみを残したパワースペクトル成分のフレーム間での相関性を評価することにより、母音の区間であるか否かを評価し、有声区間を抽出する装置であり、特徴量保存部 2 0 2 と、特徴量フレーム間相関値算出部 2 0 3 と、差分処理部 2 0 4 とを備える。なお、調波構造は、母音の発声区間内のパワースペクトル分布において主に見られる性質であり、子音の発声区間内のパワースペクトル分布においては、母音ほどの調波構造は見られない。

- 特徴量保存部 202 は、調波構造抽出部 201 より出力されるパワースペクトルを所定数のフレーム分保存する。特徴量フレーム間相関値算出部 203 は、調波構造抽出部 201 より出力されるパワースペクトルと、特徴量保存部 202 に保存されている一定フレーム前のパワースペクトルとの相関値を算出する。差分処理部 204 は、特徴量フレーム間相関値算出部 203 で求められた相関値のある一定期間における平均値を求め、特徴量フレーム間相関値算出部 203 より出力される相関値から平均値を引き、相関値と平均値との平均差分による補正相関値を求める。
- 10 音声区間決定部 205 は、差分処理部 204 より出力される平均差分による補正相関値に基づいて、音声区間を決定する。

以上のように構成された音声区間検出装置 20 の動作について以下に説明する。図 2 は、音声区間検出装置 20 が実行する処理のフローチャートである。

- 15 F F T 部 200 は、調波構造を抽出するために使用する音響特徴量として、入力信号に F F T を施すことにより、パワースペクトル成分を求める (S2)。より具体的には、F F T 部 200 は、入力信号を所定のサンプリング周波数 F_s (たとえば、11.025 kHz) でサンプリングし、1 フレーム (たとえば、10 msec) ごとに、所定のポイント
- 20 (たとえば、1 フレームあたり 128 ポイント) で F F T のスペクトル成分を求める。F F T 部 200 は、各ポイントで求められたスペクトル成分を対数化することによりパワースペクトル成分を求める。以下、パワースペクトル成分を、適宜単にスペクトル成分と表記する。

- 次に、調波構造抽出部 201 は、F F T 部 200 で抽出されたパワースペクトル成分から雑音成分等を取り除き、調波構造のみを残したパワースペクトル成分を抽出する (S4)。
- 25

FFT部200で算出されたパワースペクトル成分には、雑音によるオフセットや声道形状によって形成されるスペクトル包絡形状が含まれており、それぞれが時間変動を起こしている。このため、調波構造抽出部201は、これらの成分を取り除き、声帯振動によって形成される調波構造のみを残したパワースペクトル成分をとりだす。これにより、より効果的に有声区間検出が行なわれる。

調波構造抽出部201による処理(S4)を図3および図4を参照しながらより詳細に説明する。図3は、調波構造抽出部201による調波構造抽出処理のフローチャートであり、図4は、各フレームにおけるスペクトル成分から調波構造のみを残したスペクトル成分を抽出する過程を模式的に示す図である。

図4(a)に示されるように、調波構造抽出部201は、各フレームのスペクトル成分 $S(f)$ より、その極大値をピークホールドした値 $H_{max}(f)$ を算出し(S22)、スペクトル成分 $S(f)$ の極小値をピークホールドした値 $H_{min}(f)$ を算出する(S24)。

図4(b)に示されるように、調波構造抽出部201は、スペクトル成分 $S(f)$ から極小値のピークホールド値 $H_{min}(f)$ を引くことにより、スペクトル成分 $S(f)$ に含まれるフロア成分を除去する(S26)。これにより、雑音オフセット成分およびスペクトル包絡に起因する変動成分が除去される。

図4(c)に示されるように、調波構造抽出部201は、極大値のピークホールド値 $H_{max}(f)$ と極小値のピークホールド値 $H_{min}(f)$ との差分値を求め、ピーク変動量を算出する(S28)。

図4(d)に示されるように、調波構造抽出部201は、ピーク変動量を周波数方向に微分し、その変化量を算出する(S30)。これは、調波構造成分を有する帯域では、ピーク変動量の変化が小さいという仮定

に基づいて、調波構造の検出を行なうことを目的としている。

図 4 (e) に示されるように、調波構造抽出部 201 は、上記仮定が反映されるような重み $W(f)$ を算出する (S32)。すなわち、調波構造抽出部 201 は、ピーク変動量の変化量の絶対値と所定のしきい値とを比較し、当該変化量の絶対値が所定のしきい値 θ 以下であれば重み $W(f)$ を 1 とし、所定のしきい値 θ 以上であれば当該変化量の絶対値の逆数を重み $W(f)$ とする。これにより、ピーク変動量の変化が大きい部分の重みを小さくし、ピーク変動量の変化が小さい部分の重みを大きくすることができる。

図 4 (f) に示されるように、調波構造抽出部 201 は、フロア成分が除去されたスペクトル成分 ($S(f) - H_{min}(f)$) に重み $W(f)$ を掛け合わせ、スペクトル成分 $S'(f)$ を求める (S34)。この処理により、ピーク変動量の変化の大きい非調波構造成分を除去することが可能となる。

再度、図 2 に示される音声区間検出装置 20 の動作説明を続ける。調波構造抽出処理 (図 2 の S4、図 3) の後、特徴量フレーム間相関値算出部 203 は、調波構造抽出部 201 より出力されるスペクトル成分と、特徴量保存部 202 に保存されている所定フレーム前のスペクトル成分との間の相関値を算出する (S6)。

ここでは、着目しているフレームを j 番目のフレームとした場合、隣接するフレームのスペクトル成分を用いて相関値 $E1(j)$ を求める方法について説明する。相関値 $E1(j)$ は、次式 (1) ~ (5) に従い求められる。すなわち、 i フレームおよび $i-1$ フレームの 128 ポイントにおけるパワースペクトル成分 $P(i)$ および $P(i-1)$ を次式 (1) および (2) でそれぞれ表すものとする。また、パワースペクトル成分 $P(i)$ および $P(i-1)$ の相関関数 $xcorr(P(j-1))$,

$P(j)$ の値を次式 (3) で表すものとする。すなわち、相関関数 $xcorr(P(j-1), P(j))$ の値は、各ポイントにおける内積値からなるベクトル量である。 $z1(i)$ を次式 (4) に示されるように $xcorr(P(j-1), P(j))$ のベクトルの要素の最大値を求める。

5 これを j フレームの相関値 $E1(j)$ としてもよいし、次式 (5) で表されるようにたとえば 3 フレーム分加算した値を用いても良い。

$$P(i) = (p1(i), p2(i), \dots, p128(i)) \quad \dots (1)$$

$$P(i-1) = (p1(i-1), p2(i-1), \dots, p128(i-1)) \quad \dots (2)$$

$$xcorr(P(i-1), P(i)) =$$

$$10 \quad (p1(i-1) \times p1(i), p2(i-1) \times p2(i), \dots, p128(i-1) \times p128(i))$$

$$\dots (3)$$

$$z1(i) = \max(xcorr(P(i-1), P(i))) \quad \dots (4)$$

$$E1(j) = \sum_{i=j-2}^j z1(i) \quad \dots (5)$$

相関値 $E1(j)$ の一例を図 5 に示すグラフを用いて説明する。図 5
 15 は、入力信号を処理することにより得られる信号を表すグラフである。
 図 5 (a) は入力信号の波形を示している。この波形は、掃除機の雑音
 ($SNR = 0.5 \text{ dB}$) がある環境において、約 $1200 \sim 3000 \text{ msec}$ の間に「アールアンドビーホテルヒガシニホン」と発音している
 場合の波形である。この入力信号には、約 500 msec の箇所に掃除
 20 機を動かした際の「カタッ」という突発音が含まれ、 2800 msec
 頃に掃除機のモータの回転速度を弱から強に変更し、掃除機の音のレベ
 ルが大きくなっている。図 5 (b) は、図 5 (a) に示される入力信号

に F F T を施した場合のパワーを示しており、図 5 (c) は、相関値算出処理 (S 6) で求められた相関値の遷移を示している。

ここで、相関値 $E_1(j)$ の算出は、以下に示すような知見に基づいて算出される。すなわち、フレーム間の音響特徴量の相関値は、時間的に連続するフレームにおいて調波構造が連続していることに基づいている。このため、この調波構造を時間的に近いフレーム同士で相関をとることで、有声検出が行なわれる。調波構造が時間的に持続するのは主に母音区間である。このため、母音区間では相関値は大きくなり、子音区間では母音区間よりも相関値は小さくなるものと想定される。このように、調波構造に着目しフレーム間でパワースペクトル成分の相関値をとることによって、非周期的な雑音区間においては、相関値が小さくなるものと考えられる。このため、有声区間がより際立って識別可能となる。

また、一般的な発話スピードにおいて母音区間の持続時間は 50 ~ 150 msec (5 ~ 15 フレーム) と言われており、その持続時間内であれば、フレーム間の相関係数の値は隣接するフレームでなくとも高くなるものと想定できる。この仮定が正しければ、やはり非周期的な雑音の影響を受けにくい評価関数であるということがいえる。相関値 $E_1(j)$ を算出する際に、数フレームにわたる相関関数の値の和を用いているのは、突発的に生じる雑音の影響を除去するためと、母音であれば、上記のように 50 ~ 150 msec の持続時間があるという知見によるものである。従って、図 5 (c) に示されるように、50 フレームの近傍で発声する突発音に対しては反応せずに、相関値は小さいままである。

次に、差分処理部 204 は、特徴量フレーム間相関値算出部 203 で算出された相関値の一定時間にわたる平均値を求め、各フレームにおける相関値から当該平均値を減算し、平均差分による補正相関値を求める (S 8) 。なぜならば、相関値から平均値を引くことにより、長時間にわ

たり生じている周期性の雑音の影響を取り除くことができると考えられるためである。ここでは、5秒程度の相関値の平均値を求めており、図5(c)では、平均値を実線502で示している。すなわち、実線502よりも上の部分に相関値が存在する区間が上記平均差分による補正相関値が正の区間である。

次に、音声区間決定部205は、主に有音区間を検出する相関値E1(j)の差分処理部204で算出された平均差分による補正相関値に基づいて、後述する、相関値による選別、区間の持続長、子音区間や促音区間を加味した区間の連結、の3つの区間補正方法に従い音声区間を決定する(S10)。

ここで、音声区間決定部205による音声区間決定処理(図2のS10)についてより詳細に説明する。図6は、一発声単位で音声区間決定する処理の詳細を示すフローチャートである。

まず、第一の区間の補正方法である相関値による区間の判定について述べる。音声区間決定部205は、着目しているフレームについて、差分処理部204で求められた補正相関値が所定のしきい値よりも大きいか否かを調べる(S44)。たとえば、所定のしきい値を0とした場合には、図5(c)に示される相関値が相関値の平均値(実線502)よりも大きいか否かを調べることに等価である。

補正相関値が所定のしきい値よりも大きい場合には(S44でYES)、当該着目フレームは音声フレームであると判断し(S46)、補正相関値が所定のしきい値以下の場合には(S44でNO)、当該着目フレームは非音声フレームであると判断する(S48)。以上の音声判断処理(S44~S48)を音声区間検出対象となっているすべてのフレームについて繰返す(S42~S50)。以上の処理により、図5(d)に示されるようなグラフが得られ、音声フレームが連続する区間が有声区間として

検出される。

- このように、補正相関値の値がしきい値以下である場合には、そのフレームを非音声フレームであると判断する。ただし、騒音のレベルの影響や、音響特徴量のさまざまな条件に応じて、検出区間において期待される補正相関値が異なる。このため、音声フレームと非音声（雑音）フレームとを区別するためのしきい値は、事前の実験を通じて適宜定め用いることも可能である。この処理により調波構造性を有する信号の選別基準を厳しくすることにより、平均差分を求めた時間長より短い、例えば500ms程度の周期雑音を非音声フレームとすることが期待できる。
- 次に、第二の区間の補正方法である隣接有声区間の連結法について述べる。音声区間決定部205は、着目している有声区間と、当該有声区間に隣接する有声区間との間の距離が所定フレーム数未満であるかを調べる(S54)。たとえば、ここでは所定フレーム数を30フレームとする。当該距離が30フレーム未満の場合には(S54でYES)、隣接する2つの有声区間を連結する(S56)。以上の処理(S54~S56)をすべての有声区間について行なう(S52~S58)。以上の有声区間連結処理により、図5(e)に示されるようなグラフが得られ、近接する有声区間が連結されていることが分かる。

- 有声区間の連結をするのは、以下のような理由による。すなわち、子音区間、特に破裂音(/k/, /c/, /t/, /p/)や摩擦音などの無声子音の区間においては、調波構造が表れにくいため、相関値が小さく、有声区間として検出されにくい。しかし、子音の近傍には母音が存在するため、母音が連続する区間は有声区間とみなされるという理由による。これにより、子音部分も有声区間とすることが可能になる。

- 最後に、第三の区間の補正方法である区間持続時間について述べる。音声区間決定部205は、着目している有声区間について、その持続時

間が所定時間よりも長いかな否かを調べる（S 6 2）。たとえば、所定時間は、5 0 m s e cであるとする。持続時間が5 0 m s e cよりも長い場合には（S 6 2でY E S）、当該有声区間を音声区間と決定し（S 6 4）、持続時間が5 0 m s e c以下の場合には（S 6 2でN O）、当該有声区間
5 を非音声区間と決定する（S 6 6）。以上の処理（S 6 2～S 6 6）をすべての有声区間について行なうことにより音声区間が決定される（S 6 0～S 6 8）。以上説明した処理により、図5（f）に示すようなグラフが得られ、1 1 0～2 8 0フレームあたりに音声区間が検出される。また、図5（e）のグラフに存在していた3 2 5フレームあたりに存在し
10 ていた周期性ノイズに対する有声区間は、非音声区間と決定されていることが分かる。このように、有声区間の持続時間により有声区間を選別する処理では、相関値が高い短時間の周期的雑音を取り除くことができる。

以上説明したように本実施の形態によれば、調波構造を有するスペクトル成分のフレーム間での持続性を評価することにより、有声区間を決定している。このため、ローカルピークをトラッキングする従来の方法
15 に比べ、精度よく音声区間を決定することができる。

特に、調波構造の持続性をフレーム間のスペクトル成分の相関値により評価している。このため、フレーム間での振幅差分を取り調波構造の持続性を評価する従来方法に比べ、調波構造の有する情報を残した評価
20 が可能である。よって、短いフレームにわたる突発雑音が生じたような場合であっても、突発雑音を有声区間として検出することがない。

また、時間的に隣接する有声区間を連結することにより音声区間と決定している。このため、母音に比べ調波構造が小さい子音をも音声区間
25 と決定することが可能である。また、有声区間の持続時間を評価することにより、周期性を有する雑音を除去することが可能になる。

(実施の形態 2)

以下、図面を参照しながら本発明の実施の形態 2 に係る音声区間検出装置について説明する。本実施の形態に係る音声区間検出装置では、入力信号の S N R がよい場合には、フレーム間でのスペクトル成分の相関性のみから音声区間を決定する点が実施の形態 1 に係る音声区間検出装置とは異なる。

図 7 は、本実施の形態に係る音声区間検出装置 30 のハードウェア構成を示すブロック図である。実施の形態 1 に係る音声区間検出装置 20 と同一の構成要素については、同一の参照番号を付す。その名称および機能も同一であるため、適宜説明を省略する。なお、以下の実施の形態においても同様に適宜説明を省略する。

音声区間検出装置 30 は、入力信号の中から人間が発声している区間である音声区間を決定する装置であり、FFT 部 200 と、調波構造抽出部 201 と、有声評価部 210 と、S N R 推定部 206 と、音声区間決定部 205 とを備える。

有声評価部 210 は、有声区間を抽出する装置であり、特徴量保存部 202 と、特徴量フレーム間相関値算出部 203 と、差分処理部 204 とを備える。

S N R 推定部 206 は、差分処理部 204 より出力される平均差分による補正相関値に基づいて、入力信号の S N R を推定する。S N R 推定部 206 は、S N R が悪いと推定される場合には、差分処理部 204 より出力される補正相関値を音声区間決定部 205 に出力し、S N R がよいと推定される場合には、音声区間決定部 205 への補正相関値の出力は行わずに、差分処理部 204 より出力される補正相関値より音声区間を決定する。これは、入力信号の S N R が良好な場合には、音声区間と非音声区間との相関値の差がはっきりとしているという性質があるた

めである。

次に、SNR推定部206による入力信号のSNRの推定方法について説明する。SNR推定部206は、差分処理部204で求められる相関値の平均値が所定のしきい値未満の場合には、SNRが良好であると推定し、当該平均値が所定のしきい値以上の場合には、SNRが悪いと推定する。これは、以下のような理由に基づく。すなわち、相関値の平均値を、一発声の持続時間よりも十分に長い時間（たとえば、5秒間）にわたって求めると、SNRが良好な環境下においては、雑音区間における相関値が小さくなるため、相関値の平均値が小さくなる。これに対し、10 周期性の雑音を有するようなSNRが悪い環境下においては、雑音区間における相関値が大きくなるため、相関値の平均値が大きくなる。このように、相関値の平均値とSNRとが連動しているという性質を用いることにより、既に計算済みの一つのパラメータを評価するだけで簡単にSNRを推定することが可能である。

15 以上のように構成された音声区間検出装置30の動作について以下に説明する。図8は、音声区間検出装置30が実行する処理のフローチャートである。

FFT部200によるFFT処理(S2)から差分処理部204による補正相関値算出処理(S8)までは、図2に示した実施の形態1における音声区間検出装置20の動作と同様である。そのため、その詳細な説明はここでは繰返さない。

次に、SNR推定部206は、上記方法に従い、入力信号のSNRを推定する(S12)。SNRが良好であると推定される場合には(S14でYES)、所定のしきい値を超える補正相関値を音声区間として決定する(S16)。SNRが悪いと推定される場合には(S14でNO)、図2および図6を参照して説明した実施の形態1に係る音声区間決定部2

05による音声区間決定処理（図2のS10）と同様の処理を実行し、音声区間を決定する（S10）。

以上説明したように、本実施の形態によると、実施の形態1に記載の効果に加え、入力信号のSNRが良好な場合には、有声区間の連続性および持続時間による音声区間決定処理を行なう必要がなくなる。このため、リアルタイム性に優れた音声区間の検出が可能になる。

（実施の形態3）

以下、図面を参照しながら本発明の実施の形態3に係る音声区間検出装置について説明する。本実施の形態に係る音声区間検出装置では、調波構造性を有する音声区間を決定するのみならず、音声区間の中から特

図9は、本実施の形態に係る音声区間検出装置40のハードウェア構成を示すブロック図である。音声区間検出装置40は、入力信号の中から人間が発声している区間である音声区間と、音楽の区間である音楽区間とを決定する装置であり、FFT部200と、調波構造抽出部401と、音声・音楽区間決定部402とを備える。

調波構造抽出部401は、FFT部200で抽出されたパワースペクトル成分に基づいて、調波構造性を示す値を出力する処理部である。音声・音楽区間決定部402は、差分処理部204より出力された調波構造性を示す値に基づいて、音声区間および音楽区間を決定する処理部である。

以上のように構成された音声区間検出装置40の動作について以下に説明する。図10は、音声区間検出装置40が実行する処理のフローチャートである。

FFT部200は、調波構造を抽出するために使用する音響特徴量として、入力信号にFFTを施すことにより、パワースペクトル成分を求

める（S 2）。

次に、調波構造抽出部 4 0 1 は、FFT 部 2 0 0 で抽出されたパワースペクトル成分から、調波構造性を示す値を抽出する（S 8 2）。調波構造抽出処理（S 8 2）については、後に詳述する。

- 5 調波構造抽出部 4 0 1 は、調波構造性を示す値に基づいて、音声区間および音楽区間を決定する（S 8 4）。音声・音楽区間決定処理（S 8 4）については、後に詳述する。

- 次に、上述した調波構造抽出処理（S 8 2）について、詳細に説明する。調波構造抽出処理（S 8 2）では、パワースペクトル成分を複数の
10 帯域に分割した際に、帯域間の相関を取ることににより、調波構造性を示す値を求める。このような方法により調波構造性を示す値を求めるのは、以下のような理由による。すなわち、調波構造性は、その発生源である声帯振動における信号の影響がよく残されている帯域に見られると仮定すると、隣接帯域との間で、パワースペクトル成分の相関性が高いとい
15 う推測が成立するからである。すなわち、図 1 1 に示すように、横軸に示す各フレームにおいて、縦軸に示すパワースペクトル成分を複数の帯域（この図において、帯域数は 8）に区切った場合には、調波構造性を有する帯域間（例えば、帯域 6 0 8 と帯域 6 0 6 との間）においては、相関性が高いが、調波構造性を有しない帯域間（例えば、帯域 6 0 2 と
20 帯域 6 0 4 との間）においては、相関性が低い。

図 1 2 は、調波構造抽出処理（S 8 2）の詳細を示すフローチャートである。調波構造抽出部 4 0 1 は、各フレームについて、上述のように、各帯域間で帯域間相関値 $C(i, k)$ を算出する（S 9 2）。帯域間相関値 $C(i, k)$ は次式（6）で表される。

$$25 \quad C(i, k) = \max(Xcorr(P(i, L*(k-1)+1:L*k), P(i, L*k+1:L*(k+1))))$$

... (6)

ここで、 $P(i, x:y)$ はフレーム i のパワースペクトルにおける周波数成分 $x : y$ (x 以上、 y 以下) での、ベクトル列を示す。また、 L は帯域幅を示し、 $\max(Xcorr(\cdot))$ はベクトル列間の相関係数の最大値を示す。

- 5 調波構造型性を有する帯域では、隣接帯域との相関係数が高いため、帯域間相関値 $C(i, k)$ が大きな値を示す。逆に、調波構造型性を有しない帯域では、隣接帯域との相関係数が低いため、帯域間相関値 $C(i, k)$ が小さな値を示す。

なお、帯域間相関値 $C(i, j)$ は次式 (7) により求めてもよい。

$$10 \quad C(i, k) = \max(Xcorr(P(i, L*(k-1)+1:L*k), P(i+1, L*k+1:L*(k+1)))) \quad \dots (7)$$

- なお、式 (6) は、帯域 608 および帯域 606 間、または帯域 604 および帯域 602 間のように、同一フレーム内の隣接する帯域間でのパワースペクトルの相関を示しているのに対し、式 (7) は、帯域 608 および帯域 610 間のように、隣接するフレーム間であり、かつ隣接する帯域間でのパワースペクトルの相関を示している。式 (7) のように、隣接フレーム間でも相関を取ることにより、帯域間の相関とフレーム間の相関とを同時に計算することができる。

さらに、帯域間相関値 $C(i, k)$ は次式 (8) により求めてもよい。

$$20 \quad C(i, k) = \max(Xcorr(P(i, L*(k-1)+1:L*k), P(i, L*(k-1)+1:L*(k+1)))) \quad \dots (8)$$

式 (8) は、隣接フレームの同一帯域間でのパワースペクトルの相関を示している。

- 次に、フレーム i における調波構造型性を示す調波構造型性値 $R(i)$ および帯域番号 $N(i)$ の組 $[R(i), N(i)]$ を求める (S94)。 $[R(i), N(i)]$ は、次式 (9) に従い表される。

$$[R(i), N(i)] = [R_1(i) - R_2(i), N_1(i) - N_2(i)] \quad \dots (9)$$

ただし、 $R_1(i)$ 、 $R_2(i)$ は以下のようにあらわされる。

$$R_1(i) = \max_{k=1..L-1}(C(i,k)); \quad (10)$$

$$R_2(i) = \min_{k=1..L-1}(C(i,k)); \quad (11)$$

C : フレーム i の帯域 k における帯域調波性尺度

L : 帯域数

また、 $N_1(i)$ および $N_2(i)$ は、 $C(i, k)$ が最大となる帯域番号および最小となる帯域番号をそれぞれ示す。式 (9) に示される調波構造性値は、同一フレーム内での帯域間相関値の最大値から最小値を引くことにより求められる。このため、調波構造性のあるフレームではその値が大きくなり、調波構造性の無いフレームではその値が小さくなる。また、最大値から最小値を引くことにより、帯域間相関値を正規化している効果もある。このため、図 2 の S 8 の処理のように、平均相関値との差分処理を行なうことなく、1つのフレームにおいて正規化処理を行なうことができる。

次に、調波構造抽出部 401 は、帯域番号 $N(i)$ をその過去 X_c フレームにおける分散で重み付けした補正帯域番号 $N_d(i)$ を算出する (S 96)。また、調波構造抽出部 401 は、補正帯域番号 $N_d(i)$ の過去 X_c フレームにおける最大値 $N_e(i)$ を求める (S 98)。最大値 $N_e(i)$ を以下では重み付き帯域番号と称する。

補正帯域番号 $N_d(i)$ および重み付き帯域番号 $N_e(i)$ は $X_c = 5$ とした場合、以下の式により求められる。

$$Nd(i) = \underset{k=i-Xc_i}{\text{median}}(N(k)) - \underset{k=i-Xc_i}{\text{var}}(N(k)); \quad (12)$$

$$Ne(i) = \underset{k=i+Xc}{\text{max}}(Nd(k)); \quad (13)$$

Nd: 分散で補正した帯域番号

Ne: 分散で補正した帯域番号Ndの過去Xcフレームの最大値

5 Xc: 分散計算フレーム幅

調波構造性のない区間では、帯域番号N(i)の分散が大きくなる。

このため、補正帯域番号Nd(i)の値が小さな値(例えば、負の値)

になり、これに伴ない、重み付き帯域番号Ne(i)も小さな値になる。

さらに、調波構造抽出部401は、調波構造性値R(i)を重み付き

10 帯域番号Ne(i)で補正し、補正調波構造性値R'(i)を算出する(S100)。補正調波構造性値R'(i)は、次式(14)に従い求められる。なお、ここで用いる調波構造性値R(i)は、S8で算出した値を用いてもよい。

$$R'(i) = R(i) * Ne(i) \quad \dots (14)$$

15 図13～図15は、上述の調波構造抽出処理(S82)の実験結果を示す図である。

図13は、掃除機のノイズがある環境下(SNR=10dB)で人間が音声を発声している場合の実験結果を示す図である。40フレーム近傍には、掃除機を動かした際の「カタッ」という突発音が発生しており、

20 およそ280フレーム前後で、掃除機のモーターの回転速度を弱から強に変更したために、掃除機の音のレベルが大きくなり、周期性ノイズが発せられているものとする。また、人間は80フレームあたりから280フレームあたりまでの間に音声を発声しているものとする。

図13(a)は入力信号のパワースペクトルを示しており、図13(b) 25 は調波構造性値R(i)を示しており、図13(c)は帯域番号N(i)を示しており、図13(d)は重み付き帯域番号Ne(i)を示してお

り、図 1 3 (e) は補正調波構造性値 $R'(i)$ を示している。なお、図 1 3 (c) に示す帯域番号は、図を見やすくするために実際の帯域番号に -1 を掛けているため、0 に近いほど周波数が小さい。

図 1 3 (c) に示すように、突発音や周期性ノイズが発生している部分 (5 図中破線で囲った部分) では、帯域番号 $N(i)$ の変動が大きくなっている。このため、図 1 3 (d) に示すように、その部分の重み付き帯域番号 $N_e(i)$ は小さな値を示し、それに伴ない、図 1 3 (e) に示すように、補正調波構造性値も小さくなっている。

図 1 4 は、掃除機のノイズがほとんどない環境下 ($SNR = 40 \text{ dB}$) 10 で、図 1 3 と同じ音声を発生した場合の実験結果を示す図である。このような環境下においても図 1 3 と同様に、調波構造性のない部分の補正調波構造性値 $R'(i)$ は小さくなっている (図 1 4 (e))。

図 1 5 は、ボーカルの無い音楽に対する実験結果を示す図である。音楽では和音が出力されるため調波構造性を有するが、ドラムによりビート 15 トを刻む区間などでは調波構造性を有しない。図 1 5 (a) は入力信号のパワースペクトルを示しており、図 1 5 (b) は調波構造性値 $R(i)$ を示しており、図 1 5 (c) は帯域番号 $N(i)$ を示しており、図 1 5 (d) は重み付き帯域番号 $N_e(i)$ を示しており、図 1 5 (e) は補正調波構造性値を示している。なお、図 1 5 (c) に示す帯域番号は、20 図 1 3 (c) と同じ理由により、0 に近いほど周波数が小さい。図 1 5 (c) の破線で囲っている部分では、ドラムによりビートが刻まれることにより、調波構造性が失われている。尾のため、その部分では、図 1 5 (d) に示すように重み付き帯域番号 $N_e(i)$ が小さくなっている。したがって、図 1 5 (e) に示すように重み付き調波構造性値 $R'(i)$ 25 も小さくなっている。また、無声区間においても同様に調波構造性値 $R'(i)$ が小さくなっている。

なお、S 9 4 の処理において、フレーム i における調波構造型を示す調波構造型値 $R(i)$ および帯域番号 $N(i)$ の組 $[R(i), N(i)]$ を次式 (15) に従い求めてもよい。

$$[R(i), N(i)] = [R_1(i) - R_2(i), N_1(i) - N_2(i)] \quad \dots (15)$$

ただし、 $R_1(i)$ 、 $R_2(i)$ は以下のようにあらわされる。

$$R_1(i) = \sum_{k=1..NSP} (C(i,k)) \quad (16)$$

$$R_2(i) = \sum_{k=L-NSP..L-1} (C(i,k)) \quad (17)$$

10 C : フレーム i の帯域 k における帯域調波性尺度。

L : 帯域数

NSP : 音声ピッチ周波数帯域と仮定する帯域数

また、 $N_1(i)$ および $N_2(i)$ は、 $C(i, k)$ が最大となる帯域番号および最小となる帯域番号をそれぞれ示す。

15 なお、 $R_1(i)$ または $R_2(i)$ を調波構造型値 $R(i)$ としてもよい。

図 1 6 は、式 (15) に従い重み付き調波構造型値 $R'(i)$ を求めた実験結果である。図 1 6 は、掃除機のノイズがかなりある環境下 ($SNR = 0 \text{ dB}$) で人間が音声を発生している場合の実験結果を示す図である。なお、人間が音声を発生するタイミング、掃除機の突発音および周期性ノイズの発生タイミングは、図 1 3 に示したものと同一である。ここでは、式 (15) において、 $L = 16$ 、 $NSP = 2$ としたときの値を示している。

20 この場合においても、人間が発声しているフレームの重み付き調波構造型値 $R'(i)$ は大きい値を示し、突発音および周期性ノイズが発生しているフレームにおいては、重み付き調波構造型値 $R'(i)$ は小さい値

を示している。

次に、音声・音楽区間決定処理（図１０のＳ８４）について詳細に説明する。図１７は、音声・音楽区間決定処理（図１０のＳ８４）の詳細なフローチャートである。

5 音声・音楽区間決定部４０２は、フレーム*i*について、パワースペクトル $P(i)$ が所定の閾値 P_{min} よりも大きいかなんかを調べる（Ｓ１１２）。所定の閾値 P_{min} 以下の場合には（Ｓ１１２でＮＯ）、そのフレームは無音のフレームであると判断する（Ｓ１２６）。パワースペクトル $P(i)$ が所定の閾値 P_{min} よりも大きい場合には（Ｓ１１２でＹ

10 ＥＳ）、補正調波構造性値 $R'(i)$ が所定の閾値 R_{min} よりも大きいかなんかを判断する（Ｓ１１４）。

補正調波構造性値 $R'(i)$ が所定の閾値 R_{min} 以下の場合には（Ｓ１１４でＮＯ）、フレーム*i*が調波構造性の無い音のフレームであると判断する（Ｓ１２４）。補正調波構造性値 $R'(i)$ が所定の閾値 R_{min}

15 よりも大きい場合には（Ｓ１１４でＹＥＳ）、音声・音楽区間決定部４０２は、重み付き帯域番号 $Ne(i)$ の単位時間平均値 $ave_Ne(i)$ を算出し（Ｓ１１６）、当該単位時間平均値 $ave_Ne(i)$ が所定の閾値 Ne_min よりも大きいかなんかを調べる（Ｓ１１８）。ここで $ave_Ne(i)$ は以下の式に従い求められる。すなわち、フレーム*i*を

20 含む*d*フレーム（ここでは５０フレームとした）における $Ne(i)$ の平均値を示している。

$$ave_Ne(i) = \underset{k=i-dj}{\text{average}}(Ne(i)); \quad (18)$$

d：単位時間平均値を求めるフレーム数

25 $ave_Ne(i)$ が所定の閾値 Ne_min よりも大きい場合には（Ｓ１１８でＹＥＳ）、音楽と判断し（Ｓ１２０）、それ以外の場合には

(S 1 1 8でNO)、人間の音声のような調波構造型性を有する音であると判断する(S 1 2 2)。以上の処理(S 1 1 2～S 1 2 6)をすべてのフレームについて繰り返す(S 1 1 0～S 1 2 8)。

5 なお、以上のように $ave_Ne(i)$ の大きさにより調波構造型性を有する音の中から音楽と音声とを分離したのは以下のような考え方に基づく。すなわち、音楽も音声も信号そのものには調波構造型性を有する音であるが、音声は、有声音と無声音とが繰り返し出現される音であることより、調波構造型性値が有声音の部分では大きく、無声音の部分では小さくなり、それらが短い周期で交互に繰り返される。一方、音楽は、和
10 音が連続的に出力されるため調波構造型性を有する期間が比較的長い時間連続し、調波構造型性値が大きい状態が一定する。したがって、調波構造型性値が音楽ではあまり変動しないものの、音声では変動することを示している。換言すれば、重み付き帯域番号 $Ne(i)$ の単位時間平均値 $ave_Ne(i)$ は、音楽の方が音声よりも大きくなる。

15 なお、調波構造型性値の時間的連続性に着目して音声と音楽とを判別するようにしてもよい。すなわち、単位時間内に調波構造型性値が小さくなるフレーム数がどの程度あるかを調べるようにしてもよい。そのため、例えば、重み付き帯域番号 $Ne(i)$ が単位時間あたり負になる個数を数えるようにしてもよい。単位時間(例えば、着目しているフレーム i
20 を含む過去50フレーム)のうち、重み付き帯域番号 $Ne(i)$ が負になるフレーム数を $Ne_count(i)$ とした場合に、S 1 1 6で $ave_Ne(i)$ の代わりに $Ne_count(i)$ を算出し、S 1 1 8でフレーム数 $Ne_count(i)$ が所定の閾値よりも大きい場合に音声とし、小さい場合に音楽とするようにしてもよい。

25 以上説明したように、本実施の形態では、各フレームにおけるパワースペクトル成分を複数の帯域に区切り、帯域間で相関を取っている。こ

のため、声帯振動における信号の影響が良く残されている帯域を抽出することができ、調波構造を確実に抽出することができる。

また、調波構造の変動や、調波構造の連続性に基づいて調波構造を有する音が音楽であるのか音声であるのかを判定することができる。

5 (実施の形態 4)

次に、図面を参照しながら本発明の実施の形態 4 に係る音声区間検出装置について説明する。本実施の形態にかかる音声区間検出装置では、調波構造性値の分散に基づいて調波構造を有する音声区間を決定する。

図 18 は、本実施の形態に係る音声区間検出装置 50 のハードウェア構成を示すブロック図である。音声区間検出装置 50 は、入力信号の中から調波構造性を有する音声区間を検出する装置であり、FFT 部 200 と、調波構造抽出部 501 と、SNR 推定部 206 と、音声区間決定部 502 とを備える。

調波構造抽出部 501 は、FFT 部 200 より出力されたパワースペクトル成分に基づいて、調波構造性を示す値を出力する処理部である。音声区間決定部 502 は、調波構造性を示す値および推定された SNR に基づいて、音性区間を決定する処理部である。

以上のように構成された音声区間検出装置 50 の動作について以下に説明する。図 19 は、音声区間検出装置 50 が実行する処理のフローチャートである。FFT 部 200 は、調波構造を抽出するために使用する音響特徴量として、入力信号に FFT を施すことにより、パワースペクトル成分を求める (S2)。

次に、調波構造抽出部 501 は、FFT 部 200 で抽出されたパワースペクトル成分から、調波構造性を示す値を抽出する (S140)。調波構造処理 (S140) については、後述する。

SNR 推定部 206 は、調波構造性を示す値に基づいて、入力信号の

S N Rを推定する(S 1 2)。S N Rの推定方法は、実施の形態 2 と同様である。このため、その詳細な説明はここでは繰り返さない。

音声区間決定部 5 0 2 は、調波構造化性を示す値および推定された S N Rに基づいて音声区間を決定する(S 1 4 2)。音声区間決定処理(S 1 4 2)については、後に詳述する。

本実施の形態では、有声音と無声音との間の遷移区間に対して評価を加えることにより、音声区間決定の制度を向上させる。図 6 に示した音声区間決定方法では、(1) 音声区間間の距離が所定フレーム未満であれば、音声区間を連結し(S 5 2)、(2) 連結後の音声区間の持続時間が所定時間以下であればその区間を非音声区間としていた(S 6 0)。すなわち、無声音に対しては、(1) の処理において、S 4 2 において有声音と判断された音声の区間の間のフレームに対してなんら評価を行うことなく、(2) の処理により連結されることを暗に期待する方法である。

音声区間を詳細にみると、有声音、無声音および騒音(非音声区間)の遷移関係から次の 3 つのグループ(A グループ、B グループおよび C グループ)に分類できるものと考えられる。

A グループは有声音のグループであり、有声音から有声音への遷移、騒音から有声音への遷移、有声音から騒音への遷移が考えられる。

B グループは、有声音と無声音が混在する音のグループであり、有声音から無声音への遷移と、無声音から有声音への遷移が考えられる。

C グループは非有声音のグループであり、無声音から無声音への遷移、無声音から騒音への遷移、騒音から無声音への遷移、騒音から騒音への遷移が考えられる。

A グループに含まれる音については、調波構造化性を示す値の精度に依存して有音区間のみが決定されるものである。これに対して、B グループに含まれる音については、有音区間の周辺での音の遷移を評価するこ

とができれば、無声音区間をも抽出することが期待できるものと考えられる。C グループに含まれる音については、無声音区間だけを騒音下で抽出することは非常に難しいと考えられる。これは、騒音の性質が簡単には規定できないため、または、無声音の騒音に対する S N R が悪い場合が多いためである。

したがって、本実施の形態では、A グループのみを抽出して音声区間を決定していた図 6 の方法に加えて、有声音と無声音との間の遷移を評価することにより、B グループの音の抽出を行なう。このことにより、音声区間の決定精度を向上させることができるものとする。また、無声音から有声音への遷移区間および有声音から無声音への遷移区間において、調波構造性を示す値は大から小および小から大へとそれぞれ大きく変化していると仮定できる。このため、調波構造性を示す値を用いて有音区間と判断された区間周辺について、調波構造性を示す値の分散に基づく尺度を用いることにより、この調波構造性の値の変化を捉えることができる。ここで、調波構造性を示す値の分散を重み付き分散 V_e と呼ぶ。

次に、調波構造抽出処理（図 19 の S 140）について、詳細に説明する。図 20 は、調波構造抽出処理（S 140）の詳細を示すフローチャートである。

20 調波構造抽出部 501 は、各フレームについて、帯域間相関値 $C(i, k)$ を算出する（S 150）。帯域間相関値 $C(i, k)$ の算出は、図 12 の S 92 と同様である。このため、その詳細な説明はここでは繰り返さない。

次に、調波構造抽出部 501 は、帯域間相関値 $C(i, k)$ を用いて
25 重み付き分散 $V_e(i)$ を次式に従い算出する（S 152）。

$$Ve(i) = \text{count}(\text{if } \underset{k=1:L}{\text{var}}(\underset{j=i-Xc:i}{C(j,k)}) > th_var_change) \quad (19)$$

ここで、 Xc : フレーム幅 (= 16)

L : 帯域数 (= 16)

5 th_var_change : 閾値

である。

また、関数 $var()$ は括弧内の値の分散を示す関数であり、関数 $count()$ は、カッコ内の条件を満たす個数をカウントする関数であるものとする。

10 最後に、調波構造抽出部 501 は、調波構造型値 $R(i)$ を算出する (S154)。この算出方法は、図12のS94と同様である。このため、その詳細な説明はここでは繰り返さない。

次に、図21を参照して、音声区間決定処理 (図19のS142) について説明する。音声区間決定部 502 は、フレーム i について $R(i)$ が閾値 Th_R より大きくかつ $Ve(i)$ が閾値 Th_Ve より大きい
15 か否かを判断する (S182)。上述の条件を満たす場合には (S182でYES)、音声区間決定部 502 は、フレーム i を音声フレームであると判断し、満たさない場合には (S182でNO)、非音声フレームであると判断する (S186)。音声区間決定部 502 は、以上の処理をすべてのフレームについて行なう (S180~S188)。次に、音声区間決定部 502 は、SNR推定部 206 で推定された SNR が悪い
20 か否かを判断し (S190)、推定 SNR が悪い場合には、ループ B およびループ C の処理を実行する (S52~S68)。ループ B およびループ C の処理は図6に示したものと同様である。このため、その詳細な説明はここでは
25 は繰り返さない。

なお、推定 SNR がよい場合には (S190でNO)、ループ B を省略

し、ループCの処理（S60～S68）のみを実行する。

図22および図23は、音声区間検出装置50の実行する処理の結果を示す図である。図22は、掃除機のノイズがある環境下（SNR=10dB）で人間が音声を発声している場合の実験結果を示す図である。

5 40フレーム近傍は、掃除機を動かした際の「カタッ」という突発音が発生しており、およそ280フレーム前後で、掃除機のモーターの回転速度を弱から強に変更したために、掃除機の音のレベルが大きくなり、周期性ノイズが発せられているものとする。また、人間は80フレームあたりから280フレームあたりまでの間に音声を発声しているものとする。

図22(a)は入力信号のパワースペクトルを示しており、図22(b)は調波構造成値 $R(i)$ を示しており、図22(c)は、重み付き分散 $V_e(i)$ を示しており、図22(d)は連結前の音声区間を示しており、図22(e)は連結後の音声区間を示している。

15 図22(d)において、実線は、調波構造成値 $R(i)$ を閾値処理（図6のループA（S42～S50））することにより得られる音声区間を示しており、破線は、調波構造成値 $R(i)$ および重み付き分散 $V_e(i)$ を閾値処理（図21のループA（S180～S188））することにより得られる音声区間を示している。また、図22(e)において、破線は

20 区間連結処理（図21のS190～S68）に従い、図22(d)の破線で示した音声区間を連結した後の処理結果を示しており、実線は区間連結処理（図6のS52～S68）に従い、図22(d)の実線で示した音声区間を連結した後の処理結果を示している。図22(e)に示されるように、重み付き分散 $V_e(i)$ を用いることにより、正確に音声

25 区間を抽出することができている。

図23は、掃除機のノイズがほとんどない環境下（SNR=40dB）

で、図 2 2 と同じ音声を発生した場合の実験結果を示す図である。図 2 3 (a) ~ 図 2 3 (e) のグラフの意味は、図 2 2 (a) ~ 図 2 2 (e) のグラフの意味と同様である。図 2 3 から、区間連結前の図 2 3 (d) と区間連結後の図 2 3 (e) とを比較すると、図 2 3 (d) の破線で示される S 1 8 0 の結果は、図 2 3 (e) の実線と同様に音声区間が精度良く連結されていることを示している。したがって、推定 S N R が非常によい場合には、図 2 1 の S 1 9 0 の判定処理により、S 5 2 ~ S 5 8 の処理を行わずに、音声区間が決定されても音声区間の検出性能を維持することが可能である。

10 以上説明したように、本実施の形態によると、重み付き分散 V_e を用いて無声音と有声音との遷移区間を評価することにより、上述の B グループに属する音を抽出することができるようになった。このため、推定 S N R を用いて S N R がよいと判断された場合には区間連結を行わずとも音声区間が正確に抽出できるようになった。また、S N R が悪く、区
15 間連結が必要な場合であっても、連結時の所定フレーム数（図 2 1 の S 5 4）を小さくすることができるため、ノイズ区間を音声区間として誤検出することが少なくなった。

なお、以下に示すように調波構造型値 $R(i)$ の代わりに補正調波構造型値 $R'(i)$ を算出し、重み付き分散 $V_e(i)$ と補正調波構造型値
20 $R'(i)$ とから音声区間を検出するようにしてもよい。図 2 4 は、調波構造抽出処理（図 1 9 の S 1 4 0）の他の一例を示すフローチャートである。

調波構造抽出部 5 0 1 は、帯域間相関値 $C(i, k)$ 、重み付き分散 $V_e(i)$ および調波構造型値 $R(i)$ を算出する（S 1 6 0 ~ S 1 6 4）。
25 これらの算出方法は、図 2 0 と同様であるため、その詳細な説明はここでは繰り返さない。次に、調波構造抽出部 5 0 1 は、重み付き調波構造

性値 $R_e(i)$ を算出する (S 1 6 6)。重み付き調波構造型値 $R_e(i)$ は、次式に従い算出される。これらの式と S 9 6 / S 9 8 において算出される式との違いは、S 9 4 において算出されるフレーム i における調波構造型値 $R(i)$ を用いるかその帯域番号 $N(i)$ を用いるかの違い
 5 にある。これらの式は、ともに、重み付き分散により補正されることにより、調波構造型性を強調する指標となる。

$$R_d(i) = \text{median}_{k=i-X_c:i} (R(k)) - \text{var}_{k=i-X_c:i} (R(k)); \quad (20)$$

$$R_e(i) = \max_{k=i:i+X_c} (R_d(k)); \quad (21)$$

10 X_c : 分散計算フレーム幅 (=5)

ここで、関数 $\text{median}()$ は、括弧内の中央値を示す。

調波構造抽出部 5 0 1 は、補正調波構造型値 $R'(i)$ を算出する (S 1 6 8)。補正調波構造型値 $R'(i)$ は以下の式に従い算出される。

$$15 \quad R'(i) = R_e(i); \quad : \text{if } R_e(i) > 0; \quad (22)$$

$$R'(i) = 0; \quad : \text{if } R_e(i) < 0; \quad (23)$$

図 2 5 および図 2 6 は、図 2 4 に示したフローチャートに従い処理された処理結果を示す図である。図 2 5 は、掃除機のノイズが無い環境下
 20 (S N R = 4 0 d B) で人間が音声を発声している場合の実験結果を示しており、図 2 6 は、掃除機のノイズがある状況下 (S N R = 1 0 d B) で人間が音声を発声している場合の実験結果を示している。この実験では、図 2 3 と同じ音声を発生するものとし、突発音と周期性ノイズの発生タイミングも同じであるものとする。

25 図 2 5 (a) は入力信号を示し、図 2 5 (b) は入力信号のパワースペクトルを示しており、図 2 5 (c) は調波構造型値 $R(i)$ を示して

おり、図 25 (d) は重み付き調波構造型値 $R_e(i)$ を示しており、図 25 (e) は補正調波構造型値 $R'(i)$ を示している。図 26 (a) ~ 図 26 (e) も図 25 (a) ~ 図 25 (e) とそれぞれ同様のグラフを示している。

- 5 補正調波構造型値 $R'(i)$ は、調波構造型値 $R(i)$ 自身の分散に基づいて算出されている。このため、調波構造型を有する部分には当該分散が大きく、調波構造型を有しない部分では当該分散が小さいという性質を利用して、調波構造型を有する部分を適切に抽出することができる。
(実施の形態 5)

- 10 上述した実施の形態 1 ~ 4 に記載の音声区間決定装置では、入力信号が予めファイル等に記録されている音声に対して区間決定を行なうものである。このような処理方法は、例えば、録音済みのデータに対して処理を行なう際には、有効であるが、音声を入力しながら区間決定を行なうには不向きである。そこで、本実施の形態においては、音声の入力に
15 同期しながら音声区間をリアルタイムで決定する音声区間決定装置について説明する。

- 図 27 は、本発明の実施の形態に係る音声区間検出装置 60 の構成を示すブロック図である。音声区間検出装置 60 は、入力信号から調波構造型を有する音声区間（調波構造型区間）を検出する装置であり、FFT 部 200 と、調波構造抽出部 601 と、調波構造型区間確定部 602 と、制御部 603 とを備えている。
20 T 部 200 と、調波構造抽出部 601 と、調波構造型区間確定部 602 と、制御部 603 とを備えている。

- 図 28 は、音声区間検出装置 60 の実行する処理のフローチャートである。制御部 603 は、FR、FRS、FRE、RH、RM、CH、CM および CN を 0 にセットする (S200)。ここで、FR は、後述する
25 調波構造型値 $R(i)$ を未算出のフレームの先頭フレーム番号を示す。また、FRS は、調波構造型区間か否かが未確定の区間の先頭フレーム

番号を示す。FREは、後述する調波構造型フレーム仮判定処理を行なった最終フレームのフレーム番号を示す。RHおよびRMは調波構造型値の累積値を示す。CH、CMおよびCNはカウンタである。

FFT部200は、入力フレームをFFT変換する。調波構造抽出部
5 601は、FFT部200で抽出されたパワースペクトル成分に基づいて、調波構造型値 $R(i)$ を抽出する。以上の処理を開始フレームFRから現在時刻のフレームFRNまで行なう(S202~S210、ループA)。ループ処理が1回実行されるごとに、カウンタ i が1つずつインクリメントされ、開始フレームFRにカウンタ i の値が代入される(S
10 210)。

次に、調波構造型区間確定部602は、ここまでで求められた調波構造型値 $R(i)$ に基づいて、調波構造型を有する区間を仮判定する調波構造型フレーム仮判定処理を実行する(S212)。調波構造型フレーム仮判定処理については後述する。

15 調波構造型区間確定部602は、S212の処理の後、隣接する調波構造型区間が見つかったか否か、すなわち非調波構造型区間長CNが0より大きいかな否かを調べる(S214)。非調波構造型区間長CNは、図29(a)に図示するように、調波構造型区間の最終フレームと次の調波構造型区間の開始フレームとの間のフレーム長を示す。

20 隣接する調波構造型区間が見つかった場合には、非調波構造型区間長CNが所定の閾値よりも小さいかな否かを調べる(S216)。非調波構造型区間長CNが所定の閾値THよりも小さければ(S216でYES)、調波構造型区間確定部602は、図29(b)に示すように調波構造型区間を連結し、フレームFRS2からフレーム(FRS2+CN)まで
25 を調波構造型区間であると仮判定する(S218)。ここで、FRS2とは、非調波構造型区間であると仮判定された最初のフレーム番号を示す。

非調波構造型区間長 CN が所定の閾値 TH 以上の場合には (S 2 1 6 で NO)、図 2 9 (c) に示されるように調波構造型区間は連結されることなく、調波構造型区間確定部 6 0 2 が、後述する調波構造型区間確定処理を実行する (S 2 2 0)。その後、制御部 6 0 3 は、FSR に FRE
5 を代入し、RH、Rm、CH、CM および CN に 0 を代入する (S 2 2 2)。調波構造型区間確定処理 (S 2 2 0) については後述する。

隣接する調波構造型区間が見つからなかった場合 (S 2 1 4 で NO、図 2 9 (d))、S 2 1 8 の処理の後、または S 2 2 2 の処理の後、制御部 6 0 3 は、音声信号の入力が終了したか否かを判断する (S 2 2 4)。
10 音声信号の入力が終了していなければ (S 2 2 4 で NO)、S 2 0 2 以降の処理が繰り返される。音声信号の入力が終了していれば (S 2 2 4 で YES)、調波構造型区間確定部 6 0 2 は、調波構造型区間確定処理 (S 2 2 6) を実行し、処理を終了する。調波構造型区間確定処理 (S 2 2 6) については、後述する。

15 次に、調波構造型フレーム仮判定処理 (図 2 8 の S 2 1 2) について説明する。図 3 0 は、調波構造型フレーム仮判定処理の詳細なフローチャートである。調波構造型区間確定部 6 0 2 は、調波構造型値 $R(i)$ が予め定められた調波構造型閾値 1 よりも大きいかなんかを判断し (S 2 3 2)、大きい場合には (S 2 3 2 で YES)、着目しているフレーム i
20 を調波構造型性を有するフレームであると仮判断する。そして、累積調波構造型値 RH に調波構造型値 $R(i)$ を加算し、カウンタ CH を 1 つインクリメントする (S 2 3 4)。

次に、調波構造型区間確定部 6 0 2 は、調波構造型値 $R(i)$ が調波構造型閾値 2 よりも大きいかなんかを判断し (S 2 3 6)、大きい場合には
25 (S 2 3 6 で YES)、着目しているフレーム i を調波構造型性を有する音楽のフレームであると仮判断する。そして、累積音楽調波構造型値 RM

に調波構造型値 $R(i)$ を加算し、カウンタ CM を 1 つインクリメントする (S 2 3 6)。以上の処理をフレーム FRE からフレーム FRN まで繰り返す (S 2 3 0 ~ S 2 3 8)。

次に、調波構造型区間確定部 6 0 2 は、フレーム $FRS 2$ をフレーム FRS とした後に、着目しているフレーム i の調波構造型値 $R(i)$ が調波構造型閾値 1 よりも大きいか否かを判断し (S 2 4 2)、大きい場合にはフレーム $FRS 2$ をフレーム i とする (S 2 4 4)。以上の処理をフレーム FRS からフレーム FRN まで繰り返す (S 2 4 0 ~ S 2 4 6)。

次に、調波構造型区間確定部 6 0 2 は、カウンタ CN を 0 にセットした後に、着目しているフレーム i の調波構造型値 $R(i)$ が調波構造型閾値 1 以下であるか否かを判断し (S 2 5 0)、調波構造型閾値 1 以下である場合には (S 2 5 0 で YES)、フレーム i を非調波構造型区間であると仮判断し、カウンタ CN を 1 つインクリメントする (S 2 5 2)。以上の処理をフレーム $FRS 2$ からフレーム FRN まで繰り返す (S 2 4 8 ~ S 2 5 4)。以上の処理により、調波構造型を有する区間、音楽の調波構造型を有する区間および非調波構造型区間が仮判断される。

次に、調波構造型区間確定処理 (図 2 8 の S 2 2 0、S 2 2 6) について詳細に説明する。図 3 1 は、調波構造型区間確定処理 (図 2 8 の S 2 2 0、S 2 2 6) の詳細なフローチャートである。

調波構造型区間確定部 6 0 2 は、調波構造型を有するフレーム数を示したカウンタ CH の値が調波構造型フレーム長閾値 1 より大きく、かつ累積調波構造型値 RH が $(FRS - FRE) \times$ 調波構造型閾値 3 よりも大きいか否かを判断する (S 2 6 0)。上記条件を満たす場合には (S 2 6 0 で YES)、フレーム FRS からフレーム FRE までを調波構造型フレームであると判断する (S 2 6 2)。

調波構造型区間確定部 6 0 2 は、音楽調波構造型を有するフレーム数

を示したカウンタCMの値が調波構造型フレーム長閾値2より大きく、かつ累積音楽調波構造型値RMが $(FRS - FRE) \times$ 調波構造型閾値4よりも大きいかなかを判断する(S264)。上記条件を満たす場合には(S264でYES)、フレームFRSからフレームFREまでを音楽
5 調波構造型フレームであると判断する(S266)。

S260の条件を満たさない場合(S260でNO)、またはS264でNOの場合、音楽調波構造型は有しないが、調波構造型を有するフレームであると判断できる。このため、フレームFRSからフレームFREまでを非調波構造型フレームと判断し、カウンタCHに0を代入し、カウンタCNに $CN + FRE - FRS$ を代入する(S268)。
10

フレームワイズに調波性判断を行なう場合には調波構造型仮判定の判断を用い、より正確に調波性判断を行なう場合には調波構造型区間決定の結果を用いることにより、場合によりこれらを切り替えて使用するなどの自由度の高い選択が可能である。

15 上述したような処理を行なうことにより、調波構造型フレームと、音楽調波構造型フレームと、非調波構造型フレームと確定を行なうことができる。

以上説明したように、本実施の形態によると、入力される音声信号に対し、リアルタイムに調波構造型を有するか否かの判断を行なうことができる。このため、携帯電話などにおいて、所定フレーム遅れで非調波性のノイズを除去したりすることができる。また、音声と音楽とを見分けることができるため、携帯電話などを用いた通信において、音声部分と音楽部分とを異なる方法により符号化して通信を行なったりすることができる。
20

25 上述の実施の形態によると、環境雑音下で発声を行なった場合であっても、入力信号のレベル変動に依存せず、精度よく音声区間を決定する

ことができる。また、突発雑音や周期性雑音の影響を取り除き、精度良く音声区間を検出することができる。さらに、リアルタイムで音声区間を検出することができる。さらにまた、調波構造が小さい子音部分をも音声区間として精度良く検出することができる。また、入力信号を周波数変換したスペクトル成分にローカットフィルタをかけることにより、
5 スペクトル包絡成分を除去することができる。

以上、本発明に係る音声区間検出装置について実施の形態 1 ～ 5 に基づいて説明したが、本発明はこれらの実施の形態に限定されるものではない。

10 (FFT 部 200 の変形例)

たとえば、上述の実施の形態では、音響特徴量として FFT パワースペクトル成分を用いる方法について述べたが、FFT スペクトル成分そのものや、フレーム単位での自己相関関数や、時間軸上での線形予測残差の FFT パワースペクトル成分を用いてもよい。また、FFT スペクトルから FFT パワースペクトルを求める前に、各スペクトル成分を二乗するなどの方法により、極大値および極小値の差を拡大させ、調波構造を強調させてもよい。さらに、FFT スペクトルの対数を取り、FFT パワースペクトルを求める代わりに、FFT スペクトルの平方根を求め、FFT パワースペクトルとしてもよい。さらにまた、FFT スペクトル成分を求める前に、時間軸データに対して、フレームごとにハミング窓などの係数をかけてもよいし、プリアンファシス処理 ($1 - z^{-1}$) を行なうことで、高域強調を行ってもよい。また、音響特徴量として線スペクトル周波数 (LSF) を用いてもよい。また、周波数変換演算として、FFT に限られるものではなく、DFT (Discrete Fourier Transform)、DCT (Discrete Cosine Transform)、DST (Discrete Sine Transform) を用いても良い。
15
20
25

(調波構造抽出部 201 の変形例)

また、調波構造抽出部 201 によるスペクトル成分 $S(f)$ に含まれるフロア成分の除去処理 (図 3 の S26) の代わりに、スペクトル成分 $S(f)$ にローカットフィルタを通過させるようにしてもよい。各フレームのスペクトル成分 $S(f)$ を周波数軸方向に並べた波形とみなすと、
5 スペクトル包絡成分は、調波構造に比べゆっくりした変動である。このため、スペクトル成分にローカットフィルタをかけることにより、スペクトル包絡成分を除去することができる。この手法は時間軸上でローカットフィルタを用いて低周波数成分を取り除くことに相当するが、帯域
10 パワーやスペクトル包絡などの情報と調音構造とを同時に評価することができる点において、周波数軸上で処理する方法の方が好ましいといえる。ただし、このようなローカットフィルタを用いて算出されたスペクトル成分は、調音構造に起因する変動の他に、非周期雑音や電子音などの単一周波数を有する音声以外の音を含んでいる可能性がある。しかし、
15 これらの音は、有声評価部 210 や音声区間決定部 205 の処理により除去される。

その他のフロア成分除去の方法としては、各スペクトル成分のうち、所定の基準値以下のスペクトル成分は利用しないようにする方法がある。基準値の算出方法としては、全フレームのスペクトル成分の平均値を基準値に用いる方法、一発声の持続時間よりも十分に長い時間 (たとえば、
20 5 秒間) におけるスペクトル成分の平均値を基準値に用いる方法、スペクトル成分をいくつかの帯域に予め分割しておき、帯域ごとにスペクトル成分の平均値を求める基準値とする方法などがある。特に、静かな環境からうるさい環境へ変化するなどの環境の変動がある場合には、基準
25 値として、全フレームのスペクトル成分の平均値を利用するよりも、現在検出しようとしているフレームを含む数秒程度の区間のスペクトル成

分の平均値を用いるのがよい。

(特徴量フレーム間相関値算出部 203 の変形例)

- また、特徴量フレーム間相関値算出部 203 は、相関関数として、式 (3) の代わりに、次式 (24) を用いて相関値 $E_1(j)$ を求めるようにしてもよい。ここで、式 (24) は、 $P(i-1)$ および $P(i)$ を 128 次元ベクトル空間中のベクトルとした場合の 2 つのベクトル $P(i-1)$ および $P(i)$ がなす角の余弦を示している。また、特徴量フレーム間相関値算出部 203 は、相関値 $E_1(j)$ の代わりにフレーム j と 4 フレーム離れたフレーム間相関値を特徴とさせて、次式 (25) および (26) に従い相関値 $E_2(j)$ を求めるようにしてもよいし、8 フレーム離れたフレーム間相関値を特徴として、次式 (27) および (28) に従い相関値 $E_3(j)$ を求めるようにしてもよい。このように、離れたフレーム間で相関値を求めることにより、突発的な環境雑音に強い相関値を得ることができるという特徴がある。
- さらに、次式 (29) ~ (31) に従い、相関値 $E_1(j)$ 、相関値 $E_2(j)$ 、相関値 $E_3(j)$ の大小関係に応じた相関値 $E_4(j)$ を求めるようにしてもよいし、次式 (32) に従い相関値 $E_1(j)$ 、相関値 $E_2(j)$ 、相関値 $E_3(j)$ を加算した相関値 $E_5(j)$ を求めるようにしてもよいし、次式 (33) に従い、相関値 $E_1(j)$ 、相関値 $E_2(j)$ 、相関値 $E_3(j)$ のうちの最大値を相関値 $E_6(j)$ を求めるようにしてもよい。

$$\begin{aligned} \text{xcorr}(P(i-1), P(i)) &= \frac{P(i-1) \cdot P(i)}{|P(i-1)| |P(i)|} \\ &= \frac{p1(j-1) \times p1(j) + p2(j-1) \times p2(j) + \dots + p128(j-1) \times p128(j)}{\sqrt{p1(j-1)^2 + p2(j-1)^2 + \dots + p128(j-1)^2} \sqrt{p1(j)^2 + p2(j)^2 + \dots + p128(j)^2}} \\ &\dots (24) \end{aligned}$$

$$z2(i) = \max(\text{xcorr}(P(i-4), P(i))) \quad \dots (25)$$

$$E2(j) = \sum_{i=j-2}^j z2(i) \quad \dots (26)$$

$$z3(i) = \max(\text{xcorr}(P(i-8), P(i))) \quad \dots (27)$$

$$E3(j) = \sum_{i=j-2}^j z3(i) \quad \dots (28)$$

$$5 \quad E4(j) = z1(j) \quad \dots (29)$$

$$\text{if } (z3(j) > 0.5) \quad E4(j) = E4(j) + z1(j)/z3(j) \quad \dots (30)$$

$$\text{if } (z2(j) > 0.5) \quad E4(j) = E4(j) + z1(j)/z2(j) \quad \dots (31)$$

$$E5(j) = E1(j) + E2(j) + E3(j)$$

$$= \sum_{i=j-2}^j z1(i) + \sum_{i=j-2}^j z2(i) + \sum_{i=j-2}^j z3(i) \quad \dots (32)$$

$$10 \quad E6(j) = \max(E1(j), E2(j), E3(j))$$

$$= \max\left(\sum_{i=j-2}^j z1(i), \sum_{i=j-2}^j z2(i), \sum_{i=j-2}^j z3(i)\right) \quad \dots (33)$$

なお、相関値は、上述の $E1(j) \sim E6(j)$ の6つに限定されるわけではなく、これらの相関値を組み合わせ、新たな相関値を算出するようにしてもよい。たとえば、過去に推定された入力音響信号の SNR から、SNR が小さい場合には、相関値 $E1(j)$ を使用し、SNR が大きい場合には、相関値 $E2(j)$ または $E3(j)$ を使用するよう

15 にしてもよい。

(音声区間決定部 205 の変形例)

図 6 を用いて説明した音声区間決定部 205 の処理は、相関値による
有声区間決定処理 (S 42 ~ S 50)、有声区間の連結処理 (S 52 ~ S
58)、および有声区間の持続時間による音声区間決定処理 (S 60 ~ S
68) の 3 つの処理に大きく分類されるが、これら 3 つの処理を図 6 に
5 示される順序で実行する必要はなく、他の順序で実行するようにしても
よい。また、3 つの処理のうち、1 つまたは 2 つの処理のみを実行する
ようにしてもよい。また、図 6 は、一発声単位で処理を行なう例である
が、たとえば注目フレームごとに相関値による有声区間決定処理のみを
行なうことで、フレーム単位で音声区間を決定補正してもよい。さらに、
10 リアルタイム性が要求されることを想定して、フレーム単位の相関値に
よる音声区間を速報値として出力しておき、別途、定期的に、一発声等
長い単位で補正決定された音声区間を確定値として出力することで、リ
アルタイム性にも、検出区間性能にも対応可能な、音声検出器として作
用させてもよい。

15 (SNR 推定部 206 の変形例)

また、SNR 推定部 206 は、入力信号から直接 SNR を推定するよ
うにしてもよい。たとえば、差分処理部 204 で算出された補正相関値
が正の部分を S (シグナル) 部分とし、S 部分のパワーを求め、補正相
関値が負の部分を N (ノイズ) 部分とし、N 部分のパワーを求め、SN
20 R を求めるようにする。

(その他の変形例)

さらに、上述の音声区間検出処理を前処理とし、音声区間のみについ
て音声認識を行なう音声認識装置に音声区間検出装置を使用してもよい。

また、上述の音声区間検出処理を前処理として、音声区間のみについ
25 て録音を行なう IC (Integrated Circuit) レコーダなどの音声録音装
置に音声区間検出装置を使用しても良い。このように、音声区間のみを

録音することにより、ＩＣレコーダの記憶領域を効率的に利用することが可能となる。再生時には、音声区間のみを抽出し、話速変換機能を用いて、効率的な再生も可能となる。

また、音声区間以外の区間の入力信号をカットして雑音を抑制する雑音抑制装置に音声認識装置を利用してもよい。

さらにまた、ＶＴＲ（Video Tape Recorder）等で撮影された映像から、音声区間の映像を抽出するのに、上述の音声区間検出処理を用いてもよく、映像を編集するオーサリングツールなどにも適用可能である。

また、図４（ｆ）に示されるパワースペクトル成分 $S'(f)$ のうち、
10 調波構造が最もよく保たれている帯域を１つ以上抽出し、その帯域のみを用いて処理を行なうようにしてもよい。

また、非音声区間を検出することにより、非音声区間内でノイズの特徴を学習し、ノイズ除去のためのフィルタリング係数、ノイズ決定のパラメータ等を決めたりするようにしてもよい。このようにすることにより、
15 ノイズ除去のための装置を作成することができる。

また、上述した実施の形態における各種調波構造型値または各種相関値と、各種音声区間決定方法との組み合わせは、上述した実施の形態に限定されない。

20 産業上の利用の可能性

本発明に係る音声区間検出装置は、音声区間と雑音区間との精度よい選別が可能となるため、音声認識装置の前処理装置、音声区間のみを録音するＩＣレコーダ、音声区間と音楽区間とを異なる符号化方法で符号化する通信装置等に有用である。

請 求 の 範 囲

1. 入力音響信号から音声が含まれる区間を音声区間として検出する調波構造型音響信号区間検出方法であって、

5 前記入力音響信号に対し、所定の時間で区切られたフレーム単位で音響特徴量を抽出する音響特徴量抽出ステップと、

前記音響特徴量の持続性を評価し、評価結果に従って音声区間を決定する区間決定ステップとを含み、

10 前記音響特徴量抽出ステップでは、所定の時間で区切られたフレーム単位で前記入力音響信号を周波数変換し、調波構造を尺度化した音響特徴量を抽出し、

前記区間決定ステップでは、前記音響特徴量の同一フレーム内における相関値または前記音響特徴量の異なるフレーム間における相関値に基づいて、音声区間を決定する

15 ことを特徴とする調波構造型音響信号区間検出方法。

2. 前記音響特徴量抽出ステップでは、さらに、前記周波数変換の結果より調波構造を強調し、前記音響特徴量を抽出する

20 ことを特徴とする請求の範囲第1項に記載の調波構造型音響信号区間検出方法。

3. 前記音響特徴量抽出ステップでは、さらに、前記周波数変換の結果より調波構造を抽出し、当該調波構造を含む所定帯域の周波数変換の結果を、前記音響特徴量とする

25 ことを特徴とする請求の範囲第2項に記載の調波構造型音響信号区間検出方法。

4. 前記音響特徴量抽出ステップでは、さらに、フレーム単位の周波数変換の結果を所定の周波数帯域幅ごとに分割し、同一フレーム内の所定の周波数帯域間で、前記周波数変換の結果の相関値を算出し、算出結果に基づいて前記音響特徴量を抽出する

5 ことを特徴とする請求の範囲第1項に記載の調波構造型音響信号区間検出方法。

5. 前記音響特徴量抽出ステップでは、さらに、フレームごとに前記相関値の最大値と最小値との差を求め、当該差に基づいて、前記音響特徴量を抽出する

10 ことを特徴とする請求の範囲第4項に記載の調波構造型音響信号区間検出方法。

6. 前記音響特徴量抽出ステップでは、フレーム単位の周波数変換の結果を所定の周波数帯域幅ごとに分割し、異なるフレーム間、かつ所定の周波数帯域間で、前記周波数変換の結果の相関値を算出し、算出結果に基づいて前記音響特徴量を抽出する

20 ことを特徴とする請求の範囲第1項に記載の調波構造型音響信号区間検出方法。

7. 前記音響特徴量抽出ステップでは、さらに、フレームごとに前記相関値の最大値と最小値との差を求め、前記差に基づいて、前記音響特徴量を抽出する

25 ことを特徴とする請求の範囲第6項に記載の調波構造型音響信号区間検出方法。

8. 前記区間決定ステップでは、異なるフレーム間における前記音響特徴量の相関値に基づいて、前記音響特徴量の持続性を評価し、評価結果に従って音声区間を決定する

5 ことを特徴とする請求の範囲第1項に記載の調波構造型音響信号区間検出方法。

9. 前記区間決定ステップでは、異なるフレーム間における前記音響特徴量の分散に基づいて、前記音響特徴量の持続性を評価し、評価結果に従って音声区間を決定する

10 ことを特徴とする請求の範囲第1項に記載の調波構造型音響信号区間検出方法。

10. 前記音響特徴量の持続性を評価する評価値を算出する評価ステップと、

15 前記評価値の時間的な連続性を評価し、評価結果に従って音声区間を決定する音声区間決定ステップとを含む

ことを特徴とする請求の範囲第1項に記載の調波構造型音響信号区間検出方法。

20 11. 前記区間決定ステップは、さらに、

所定数のフレームにわたる、音響特徴量抽出ステップにおいて算出される音響特徴量または、前記評価ステップにおいて算出される前記評価値と、第1の所定しきい値との比較に基づいて、前記入力音響信号の音声雑音比を推定するステップと、

25 推定された前記音声雑音比が第2の所定しきい値以上の場合には、前記評価ステップにおいて算出される前記評価値に基づいて前記音声区間

を決定するステップとを含み、

前記音声区間決定ステップでは、前記音声雑音比が前記第 2 の所定しきい値未満の場合に、前記評価値の時間的な連続性を評価し、評価結果に従って前記音声区間を決定する

- 5 ことを特徴とする請求の範囲第 10 項に記載の調波構造型音響信号区間検出方法。

12. 前記区間決定ステップは、

前記音響特徴量の持続性を評価する評価値を算出する評価ステップと、

- 10 前記評価値の時間的な連続性を評価し、評価結果に従って調波構造を有するが音声ではない非音声調波構造区間を決定する非音声調波構造区間決定ステップとを含む

ことを特徴とする請求の範囲第 1 項に記載の調波構造型音響調波構造型音響信号区間検出方法。

15

13. 前記音響特徴量抽出ステップは、

所定の時間で区切られたフレーム単位で、前記入力音響信号を周波数変換する周波数変換ステップと、

- 20 フレーム単位の周波数変換の結果を所定の周波数帯域ごとに分割し、同一フレーム内の所定の周波数帯域間で、前記周波数変換の結果の相関値を算出する相関値算出ステップと、

同一フレーム内における相関値の最大値または最小値をとる周波数帯域の識別子を前記音響特徴量として抽出する抽出ステップとを含む

- 25 ことを特徴とする請求の範囲第 12 項に記載の調波構造型音響信号区間検出方法。

14. 前記音響特徴量抽出ステップは、

所定の時間で区切られたフレーム単位で、前記入力音響信号を周波数変換する周波数変換ステップと、

5 所定数離れたフレーム間で前記周波数変換の結果の相関値を算出する相関値算出ステップと、

所定数のフレームごとに前記相関値の分散を算出することにより、前記調波構造を尺度化した音響特徴量を抽出する音響特徴量抽出ステップを含む

10 ことを特徴とする請求の範囲第1項に記載の調波構造音響区間検出方法。

15. 前記区間決定ステップでは、2種以上の異なる間隔のフレーム間における複数の相関値に基づいて、前記持続性を評価する

15 ことを特徴とする請求の範囲第1項に記載の調波構造型音響信号区間検出方法。

16. 前記区間決定ステップでは、前記入力音響信号の音声雑音比に基づいて、前記2種以上の異なる間隔のフレーム間における相関値のいずれかを選択し、選択された相関値に基づいて前記持続性を評価する

20 ことを特徴とする請求の範囲第15項に記載の調波構造型音響信号区間検出方法。

17. 前記区間決定ステップでは、前記音響特徴量のフレーム間における相関値と、前記相関値を所定フレーム数にわたり平均した平均値との補正相関値に基づいて、前記持続性を評価する

ことを特徴とする請求の範囲第1項に記載の調波構造型音響信号区間

検出方法。

18. 前記音響特徴量抽出手段は、所定の時間で区切られたフレーム単位で前記入力音響信号を周波数変換し、調波構造を尺度化した音響特徴量を抽出し、

前記区間決定手段は、前記音響特徴量の同一フレーム内における相関値または前記音響特徴量の異なるフレーム間における相関値に基づいて、音声区間を決定する

10 ことを特徴とする請求の範囲第28項に記載の調波構造型音響信号区間検出装置。

19. 入力音響信号に含まれる音声を認識する音声認識装置であって、

15 前記入力音響信号に対し、所定の時間で区切られたフレーム単位で音響特徴量を抽出する音響特徴量抽出手段と、

前記音響特徴量の持続性を評価し、評価結果に従って音声区間を決定する区間決定手段と、

前記区間決定手段で決定された音声区間において音声認識を行なう認識手段とを備え、

20 前記音響特徴量抽出手段は、所定の時間で区切られたフレーム単位で前記入力音響信号を周波数変換し、調波構造を尺度化した音響特徴量を抽出し、

前記区間決定手段は、前記音響特徴量の同一フレーム内における相関値または前記音響特徴量の異なるフレーム間における相関値に基づいて、

25 音声区間を決定する

ことを特徴とする音声認識装置。

20. 入力音響信号に含まれる音声を録音する音声録音装置であって、

前記入力音響信号に対し、所定の時間で区切られたフレーム単位で音響特徴量を抽出する音響特徴量抽出手段と、

5 前記音響特徴量の持続性を評価し、評価結果に従って音声区間を決定する区間決定手段と、

前記区間決定手段で決定された音声区間における入力音響信号を録音する録音手段とを備え、

10 前記音響特徴量抽出手段は、所定の時間で区切られたフレーム単位で前記入力音響信号を周波数変換し、調波構造を尺度化した音響特徴量を抽出し、

前記区間決定手段は、前記音響特徴量の同一フレーム内における相関値または前記音響特徴量の異なるフレーム間における相関値に基づいて、音声区間を決定する

15 ことを特徴とする音声録音装置。

21. 入力音響信号に対し、所定の時間で区切られたフレーム単位で音響特徴量を抽出する音響特徴量抽出ステップと、

20 前記音響特徴量の持続性を評価し、評価結果に従って音声区間を決定する区間決定ステップとをコンピュータに実行させ、

前記音響特徴量抽出ステップでは、所定の時間で区切られたフレーム単位で前記入力音響信号を周波数変換し、調波構造を尺度化した音響特徴量を抽出し、

25 前記区間決定ステップでは、前記音響特徴量の同一フレーム内における相関値または前記音響特徴量の異なるフレーム間における相関値に基づいて、音声区間を決定する

ことを特徴とするプログラム。

補正書の請求の範囲

[2004年11月1日(01.11.04)国際事務局受理:出願当初の請求の範囲
18は補正された;他の請求の範囲は変更なし。(8頁)]

1. 入力音響信号から音声が含まれる区間を音声区間として検出する調波構造型音響信号区間検出方法であって、

5 前記入力音響信号に対し、所定の時間で区切られたフレーム単位で音響特徴量を抽出する音響特徴量抽出ステップと、

前記音響特徴量の持続性を評価し、評価結果に従って音声区間を決定する区間決定ステップとを含み、

10 前記音響特徴量抽出ステップでは、所定の時間で区切られたフレーム単位で前記入力音響信号を周波数変換し、調波構造を尺度化した音響特徴量を抽出し、

前記区間決定ステップでは、前記音響特徴量の同一フレーム内における相関値または前記音響特徴量の異なるフレーム間における相関値に基づいて、音声区間を決定する

15 ことを特徴とする調波構造型音響信号区間検出方法。

2. 前記音響特徴量抽出ステップでは、さらに、前記周波数変換の結果より調波構造を強調し、前記音響特徴量を抽出する

20 ことを特徴とする請求の範囲第1項に記載の調波構造型音響信号区間検出方法。

3. 前記音響特徴量抽出ステップでは、さらに、前記周波数変換の結果より調波構造を抽出し、当該調波構造を含む所定帯域の周波数変換の結果を、前記音響特徴量とする

25 ことを特徴とする請求の範囲第2項に記載の調波構造型音響信号区間検出方法。

4. 前記音響特徴量抽出ステップでは、さらに、フレーム単位の周波数変換の結果を所定の周波数帯域幅ごとに分割し、同一フレーム内の所定の周波数帯域間で、前記周波数変換の結果の相関値を算出し、算出結果に基づいて前記音響特徴量を抽出する

5 ことを特徴とする請求の範囲第1項に記載の調波構造的音響信号区間検出方法。

5. 前記音響特徴量抽出ステップでは、さらに、フレームごとに前記相関値の最大値と最小値との差を求め、当該差に基づいて、前記音響特徴量を抽出する

ことを特徴とする請求の範囲第4項に記載の調波構造的音響信号区間検出方法。

6. 前記音響特徴量抽出ステップでは、フレーム単位の周波数変換の結果を所定の周波数帯域幅ごとに分割し、異なるフレーム間、かつ所定の周波数帯域間で、前記周波数変換の結果の相関値を算出し、算出結果に基づいて前記音響特徴量を抽出する

ことを特徴とする請求の範囲第1項に記載の調波構造的音響信号区間検出方法。

7. 前記音響特徴量抽出ステップでは、さらに、フレームごとに前記相関値の最大値と最小値との差を求め、前記差に基づいて、前記音響特徴量を抽出する

ことを特徴とする請求の範囲第6項に記載の調波構造的音響信号区間検出方法。

8. 前記区間決定ステップでは、異なるフレーム間における前記音響特徴量の相関値に基づいて、前記音響特徴量の持続性を評価し、評価結果に従って音声区間を決定する

5 ことを特徴とする請求の範囲第1項に記載の調波構造型音響信号区間検出方法。

9. 前記区間決定ステップでは、異なるフレーム間における前記音響特徴量の分散に基づいて、前記音響特徴量の持続性を評価し、評価結果に従って音声区間を決定する

10 ことを特徴とする請求の範囲第1項に記載の調波構造型音響信号区間検出方法。

10. 前記音響特徴量の持続性を評価する評価値を算出する評価ステップと、

15 前記評価値の時間的な連続性を評価し、評価結果に従って音声区間を決定する音声区間決定ステップとを含む

ことを特徴とする請求の範囲第1項に記載の調波構造型音響信号区間検出方法。

20 11. 前記区間決定ステップは、さらに、

所定数のフレームにわたる、音響特徴量抽出ステップにおいて算出される音響特徴量または、前記評価ステップにおいて算出される前記評価値と、第1の所定しきい値との比較に基づいて、前記入力音響信号の音声雑音比を推定するステップと、

25 推定された前記音声雑音比が第2の所定しきい値以上の場合には、前記評価ステップにおいて算出される前記評価値に基づいて前記音声区間

を決定するステップとを含み、

前記音声区間決定ステップでは、前記音声雑音比が前記第2の所定しきい値未満の場合に、前記評価値の時間的な連続性を評価し、評価結果に従って前記音声区間を決定する

- 5 ことを特徴とする請求の範囲第10項に記載の調波構造型音響信号区間検出方法。

12. 前記区間決定ステップは、

- 前記音響特徴量の持続性を評価する評価値を算出する評価ステップと、
10 前記評価値の時間的な連続性を評価し、評価結果に従って調波構造を有するが音声ではない非音声調波構造区間を決定する非音声調波構造区間決定ステップとを含む

ことを特徴とする請求の範囲第1項に記載の調波構造型音響調波構造型音響信号区間検出方法。

15

13. 前記音響特徴量抽出ステップは、

所定の時間で区切られたフレーム単位で、前記入力音響信号を周波数変換する周波数変換ステップと、

- フレーム単位の周波数変換の結果を所定の周波数帯域ごとに分割し、
20 同一フレーム内の所定の周波数帯域間で、前記周波数変換の結果の相関値を算出する相関値算出ステップと、

同一フレーム内における相関値の最大値または最小値をとる周波数帯域の識別子を前記音響特徴量として抽出する抽出ステップとを含む

- ことを特徴とする請求の範囲第12項に記載の調波構造型音響信号区
25 間検出方法。

14. 前記音響特徴量抽出ステップは、

所定の時間で区切られたフレーム単位で、前記入力音響信号を周波数変換する周波数変換ステップと、

5 所定数離れたフレーム間で前記周波数変換の結果の相関値を算出する相関値算出ステップと、

所定数のフレームごとに前記相関値の分散を算出することにより、前記調波構造を尺度化した音響特徴量を抽出する音響特徴量抽出ステップとを含む

10 ことを特徴とする請求の範囲第1項に記載の調波構造音響区間検出方法。

15. 前記区間決定ステップでは、2種以上の異なる間隔のフレーム間における複数の相関値に基づいて、前記持続性を評価する

15 ことを特徴とする請求の範囲第1項に記載の調波構造型音響信号区間検出方法。

16. 前記区間決定ステップでは、前記入力音響信号の音声雑音比に基づいて、前記2種以上の異なる間隔のフレーム間における相関値のいずれかを選択し、選択された相関値に基づいて前記持続性を評価する

20 ことを特徴とする請求の範囲第15項に記載の調波構造型音響信号区間検出方法。

17. 前記区間決定ステップでは、前記音響特徴量のフレーム間における相関値と、前記相関値を所定フレーム数にわたり平均した平均値
25 との補正相関値に基づいて、前記持続性を評価する

ことを特徴とする請求の範囲第1項に記載の調波構造型音響信号区間

検出方法。

18. (補正後) 入力音響信号から音声が含まれる区間を音声区間として検出する調波構造型音響信号区間検出装置であって、

5 前記入力音響信号に対し、所定の時間で区切られたフレーム単位で音響特徴量を抽出する音響特徴量抽出手段と、

前記音響特徴量の持続性を評価し、評価結果に従って音声区間を決定する区間決定手段とを含み、

10 前記音響特徴量抽出手段は、所定の時間で区切られたフレーム単位で前記入力音響信号を周波数変換し、調波構造を尺度化した音響特徴量を抽出し、

前記区間決定手段は、前記音響特徴量の同一フレーム内における相関値または前記音響特徴量の異なるフレーム間における相関値に基づいて、音声区間を決定する

15 ことを特徴とする調波構造型音響信号区間検出装置。

19. 入力音響信号に含まれる音声を認識する音声認識装置であって、

20 前記入力音響信号に対し、所定の時間で区切られたフレーム単位で音響特徴量を抽出する音響特徴量抽出手段と、

前記音響特徴量の持続性を評価し、評価結果に従って音声区間を決定する区間決定手段と、

前記区間決定手段で決定された音声区間において音声認識を行なう認識手段とを備え、

25 前記音響特徴量抽出手段は、所定の時間で区切られたフレーム単位で前記入力音響信号を周波数変換し、調波構造を尺度化した音響特徴量を

抽出し、

前記区間決定手段は、前記音響特徴量の同一フレーム内における相関値または前記音響特徴量の異なるフレーム間における相関値に基づいて、音声区間を決定する

5 ことを特徴とする音声認識装置。

20. 入力音響信号に含まれる音声を録音する音声録音装置であって、

10 前記入力音響信号に対し、所定の時間で区切られたフレーム単位で音響特徴量を抽出する音響特徴量抽出手段と、

前記音響特徴量の持続性を評価し、評価結果に従って音声区間を決定する区間決定手段と、

前記区間決定手段で決定された音声区間における入力音響信号を録音する録音手段とを備え、

15 前記音響特徴量抽出手段は、所定の時間で区切られたフレーム単位で前記入力音響信号を周波数変換し、調波構造を尺度化した音響特徴量を抽出し、

20 前記区間決定手段は、前記音響特徴量の同一フレーム内における相関値または前記音響特徴量の異なるフレーム間における相関値に基づいて、音声区間を決定する

ことを特徴とする音声録音装置。

21. 入力音響信号に対し、所定の時間で区切られたフレーム単位で音響特徴量を抽出する音響特徴量抽出ステップと、

25 前記音響特徴量の持続性を評価し、評価結果に従って音声区間を決定する区間決定ステップとをコンピュータに実行させ、

前記音響特徴量抽出ステップでは、所定の時間で区切られたフレーム単位で前記入力音響信号を周波数変換し、調波構造を尺度化した音響特徴量を抽出し、

- 5 前記区間決定ステップでは、前記音響特徴量の同一フレーム内における相関値または前記音響特徴量の異なるフレーム間における相関値に基づいて、音声区間を決定する

ことを特徴とするプログラム。

条約 19 条に基づく説明書

補正後の請求の範囲第 18 項は、補正前の請求の範囲第 1 項に記載の調波構造型音響信号区間検出方法に対応する調波構造型音響信号区間検出装置の請求の範囲である。

補正後の請求の範囲第 18 項は、補正前の請求の範囲第 1 項に記載の調波構造型音響信号区間検出方法に含まれる「ステップ」の文言を、「手段」に置き換えたものである。したがって、請求の範囲第 18 項に対する補正は、新規事項の追加ではない。

図1

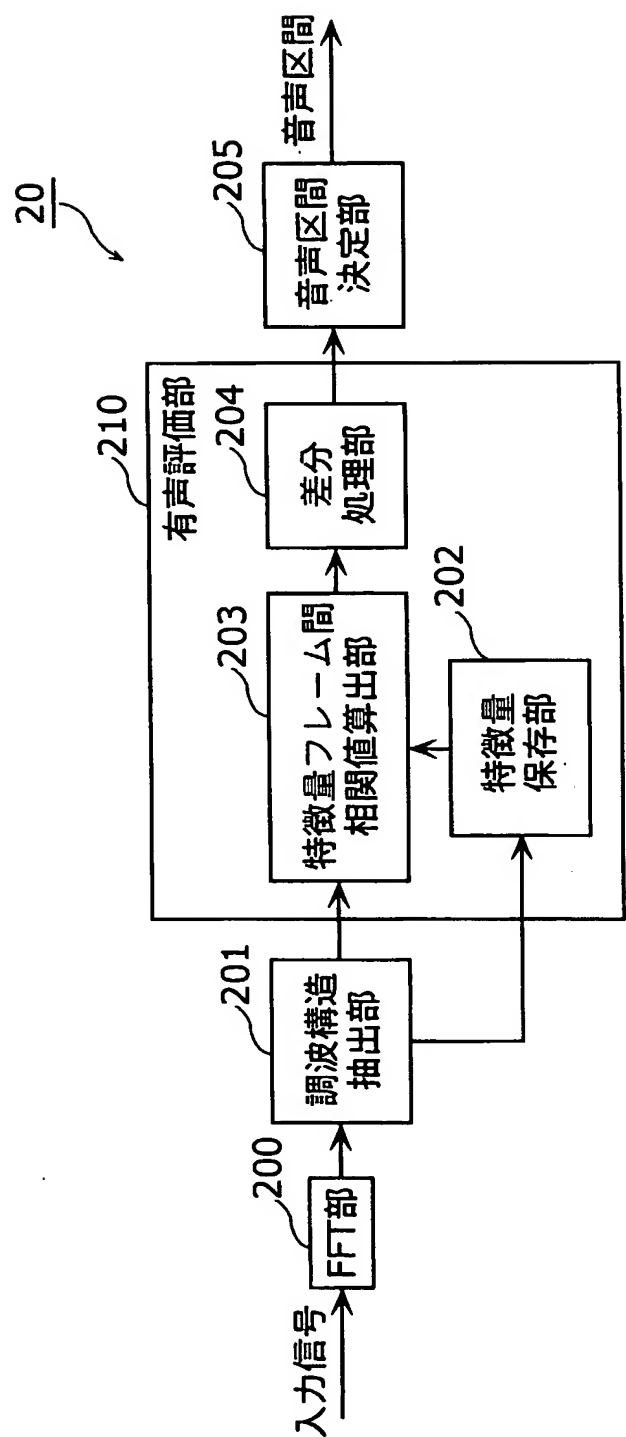


図2

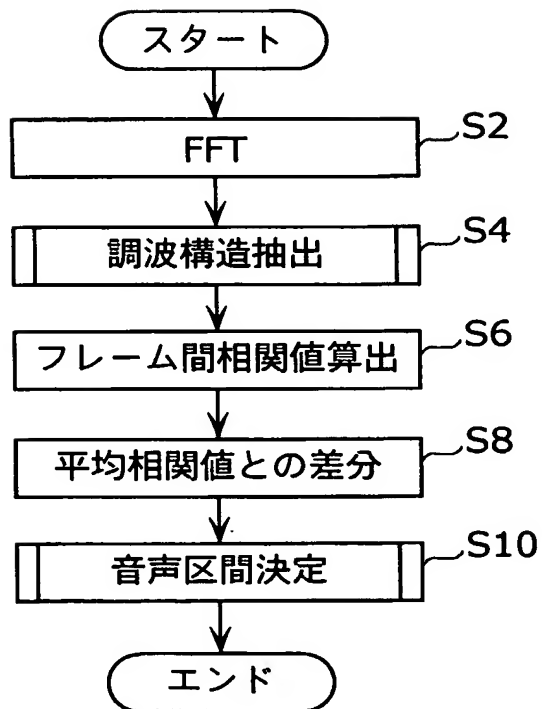


図3

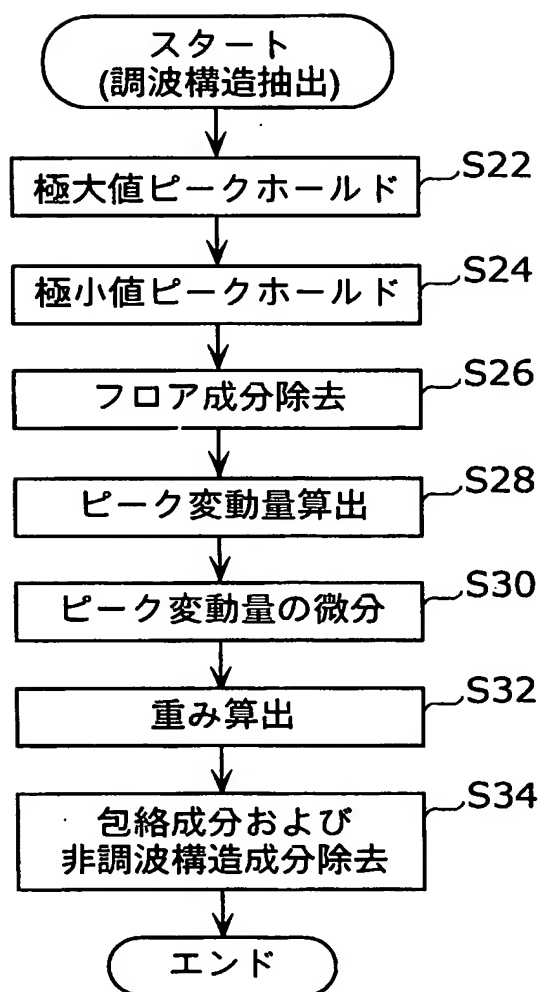


図4

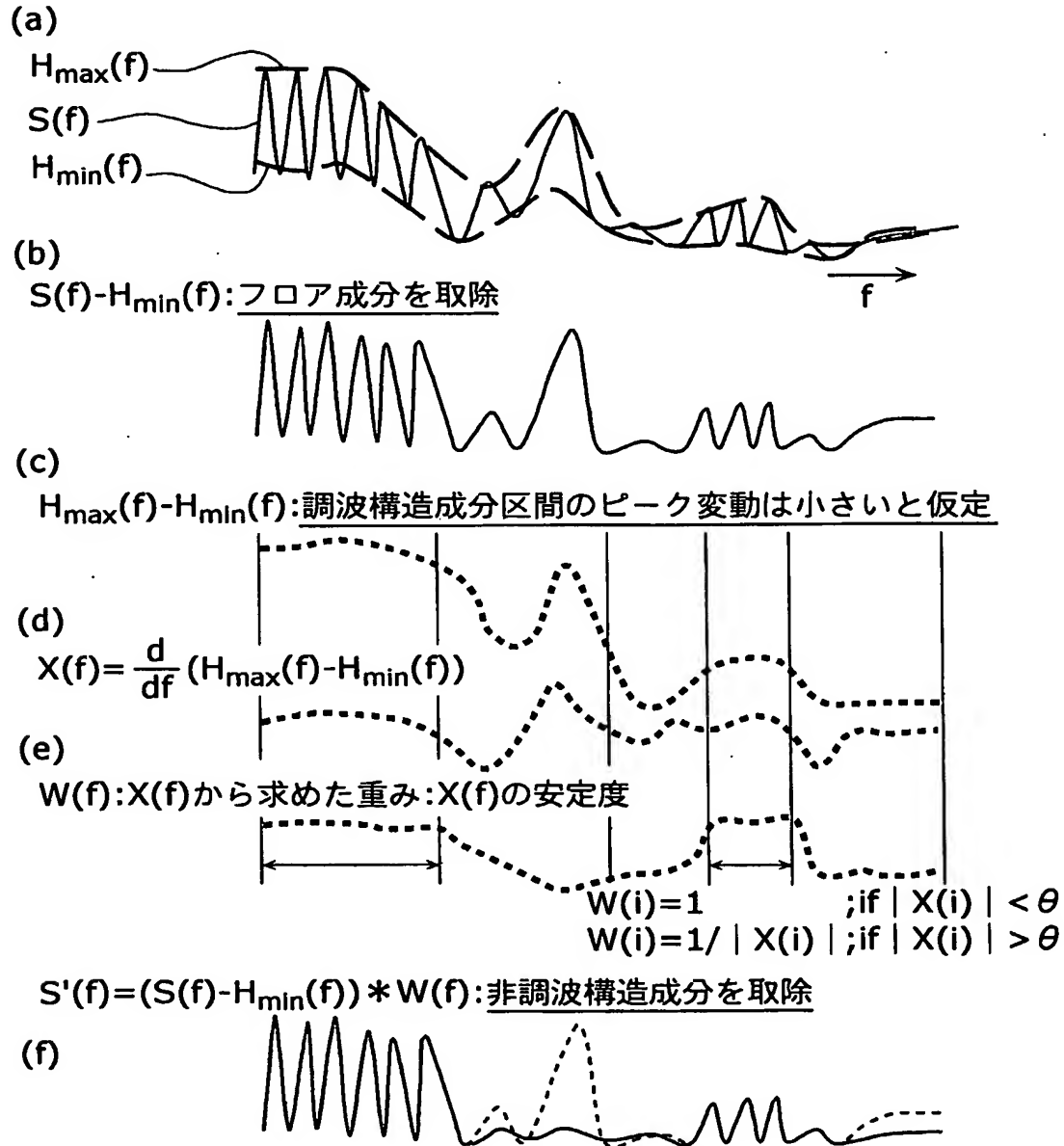


図5

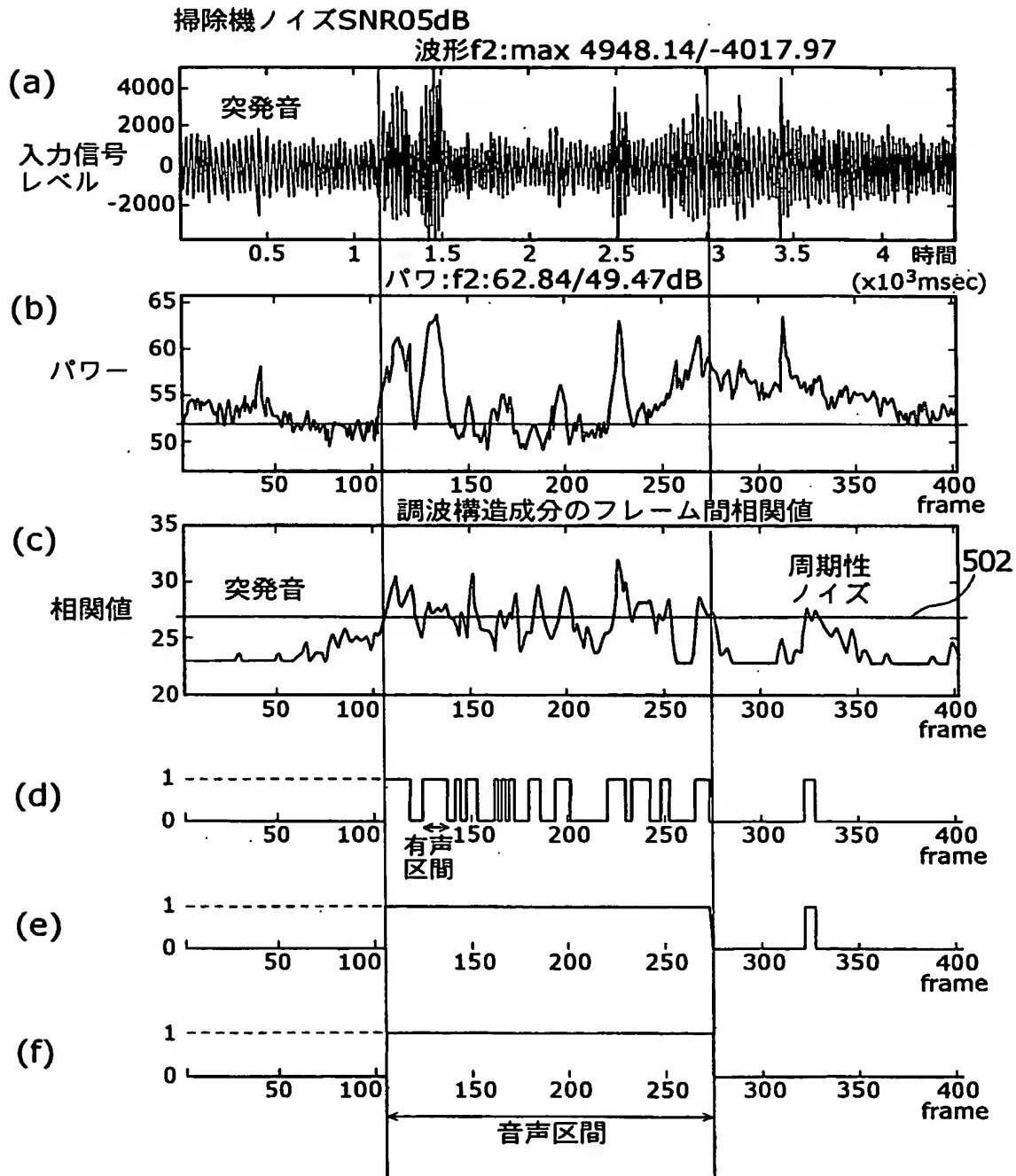


図6

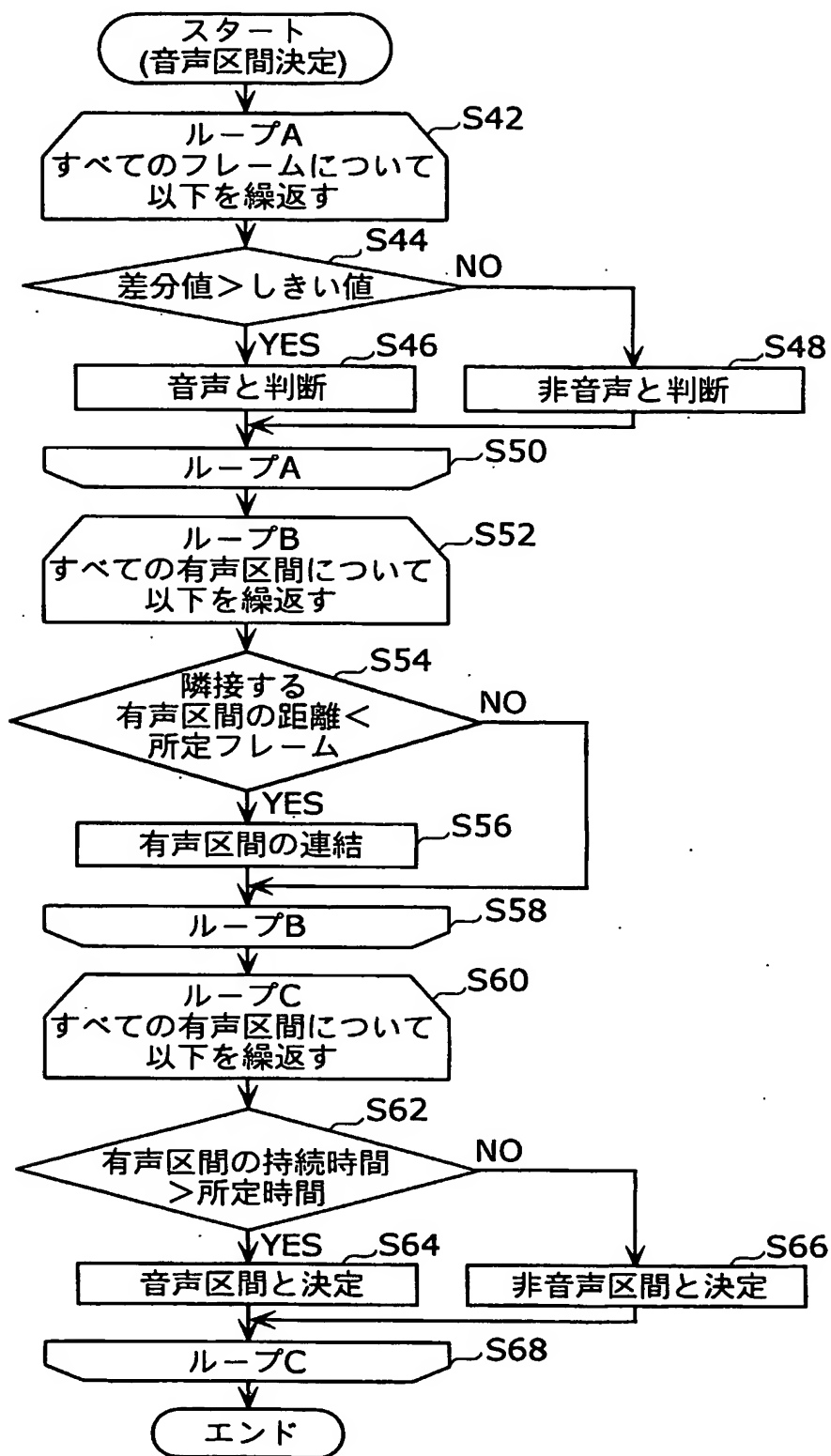


図7

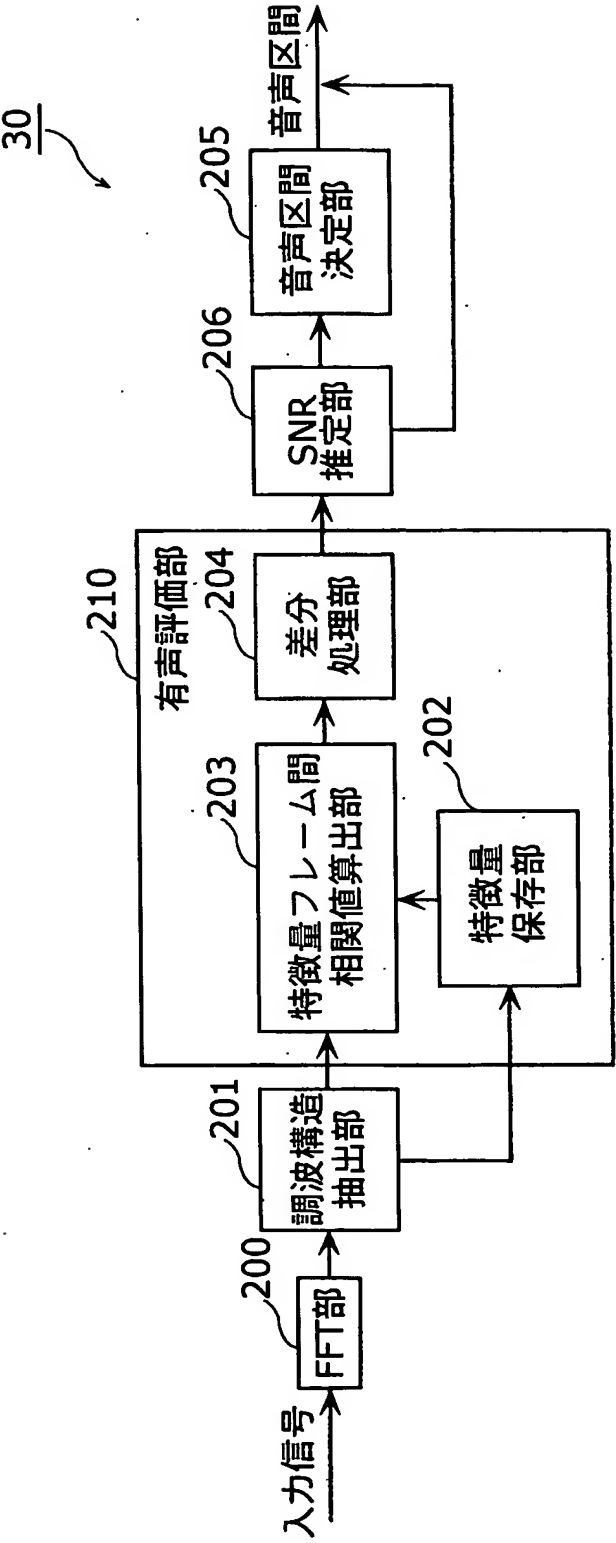


図8

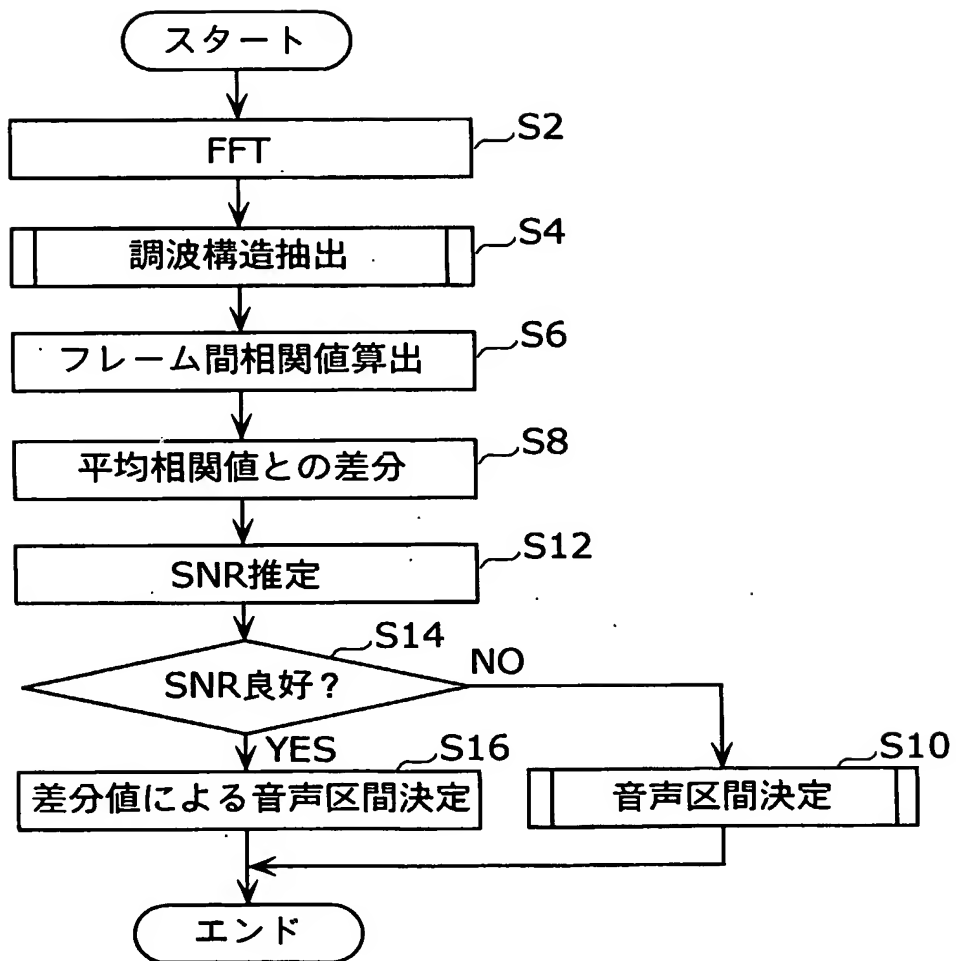


図9

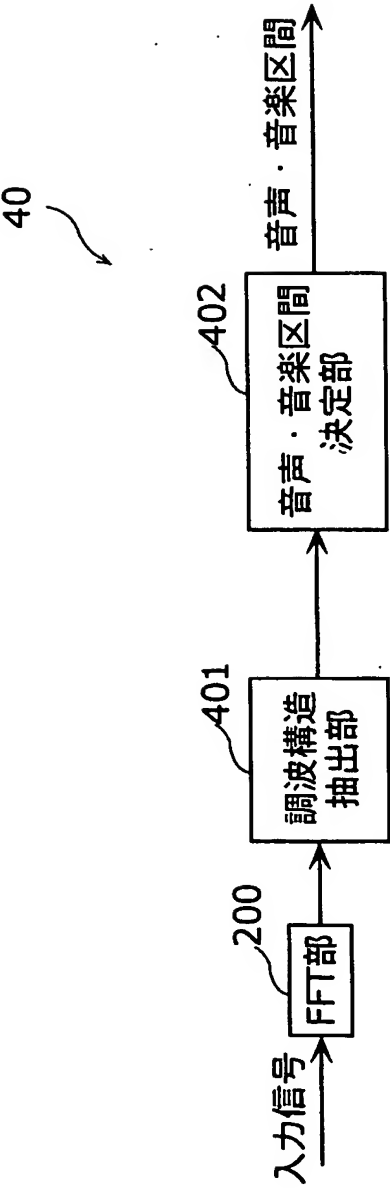


図10

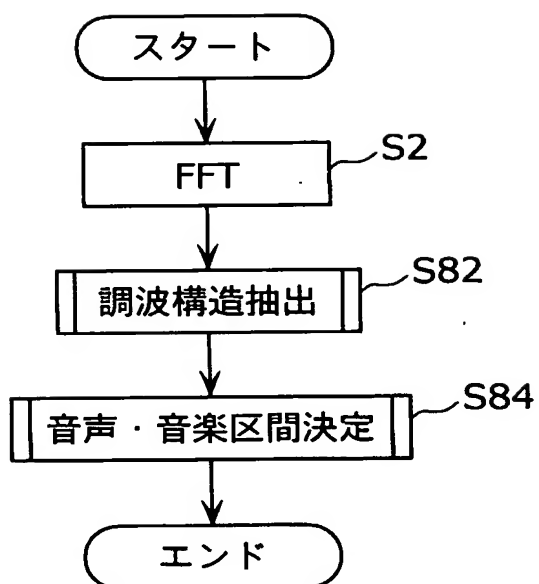


図11

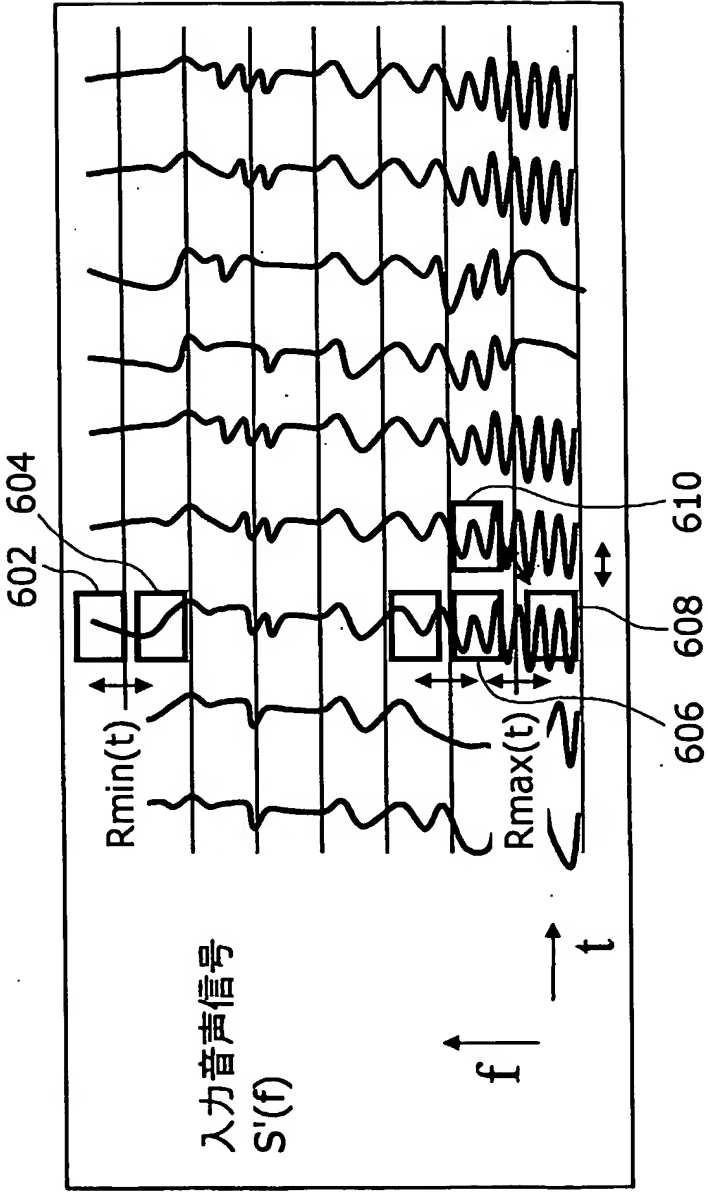


図12

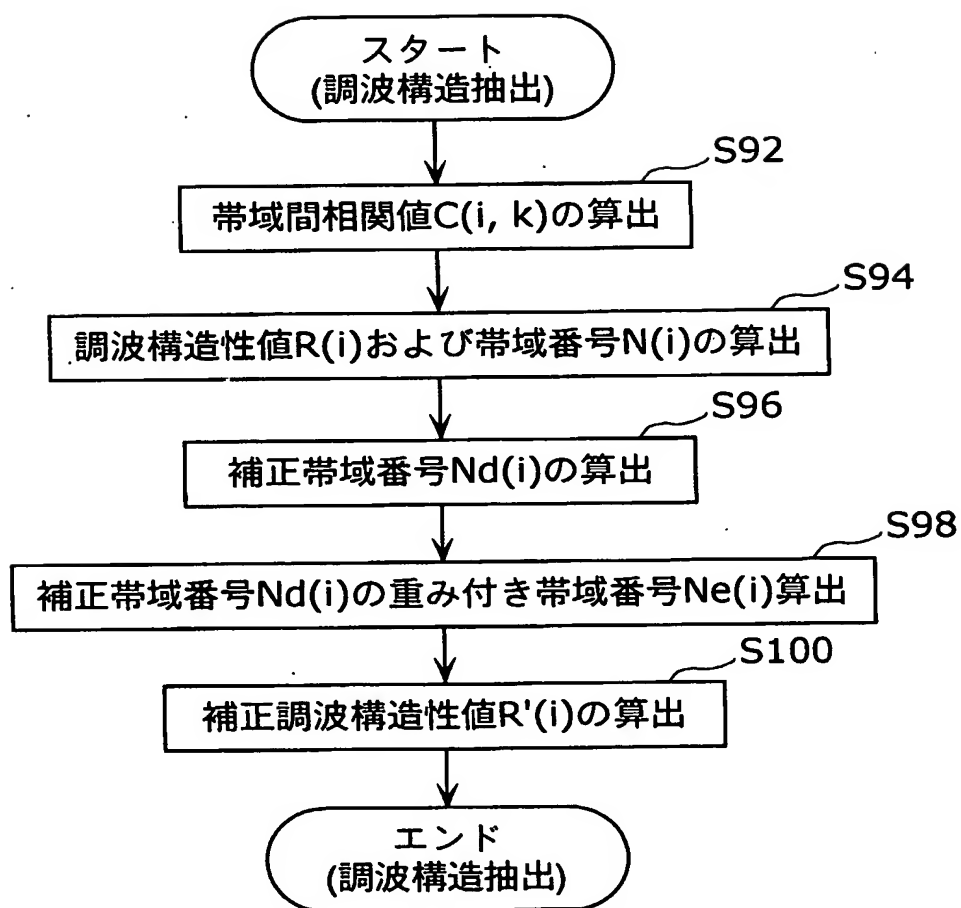


図13

音声: 掃除機ノイズ10dB下(ノイズ有)

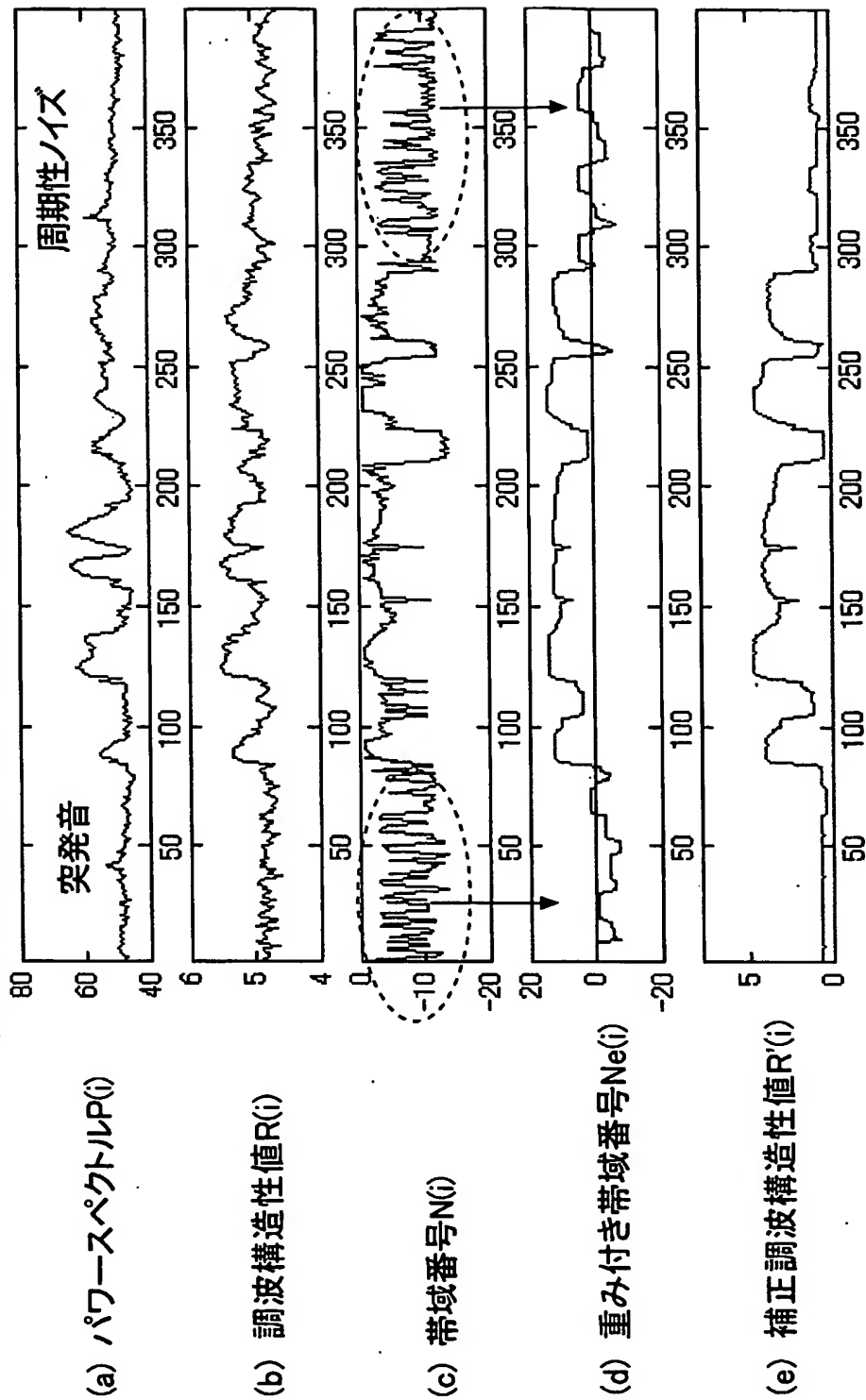


図14

音声1: 掃除機ノイズ40dB下(ノイズ無)

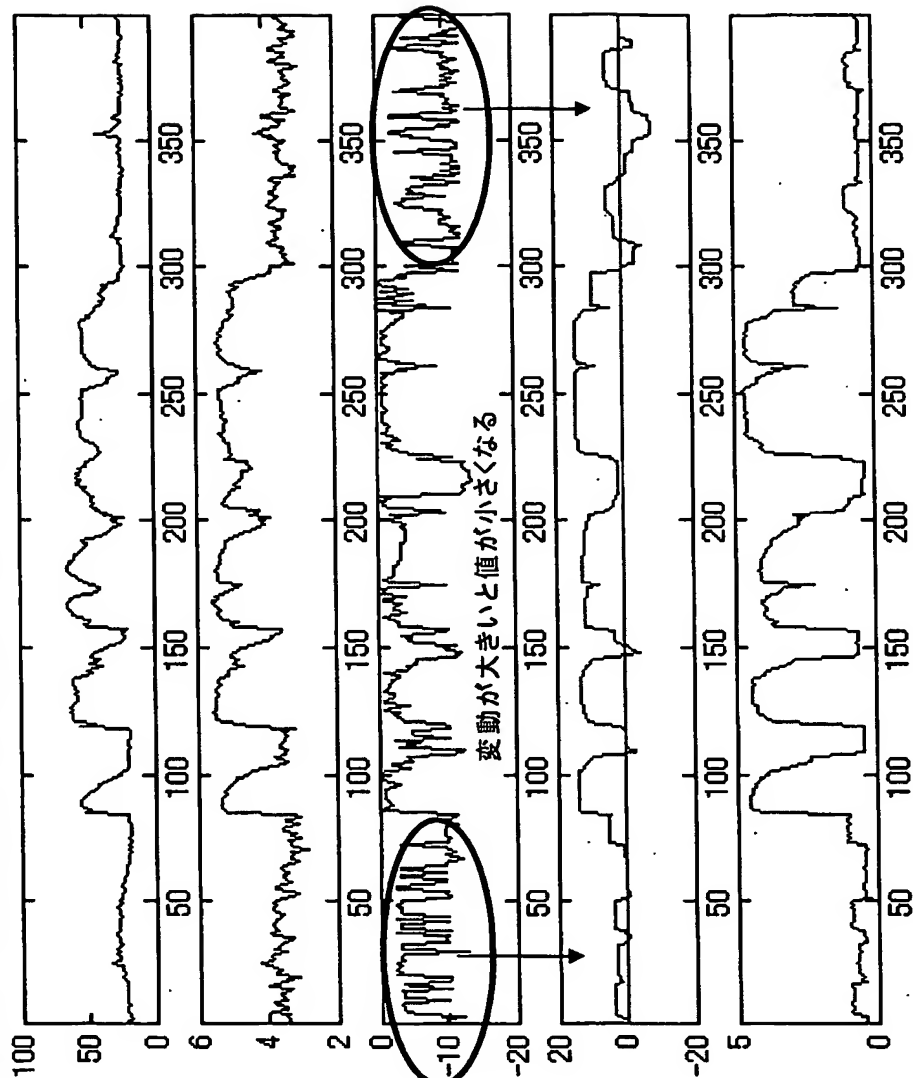
(a) パワースペクトル $P(i)$ (b) 調波構造性値 $R(i)$ (c) 帯域番号 $Ne(i)$ (d) 重み付き帯域番号 $Ne(i)$ (e) 補正調波構造性値 $R'(i)$

図15

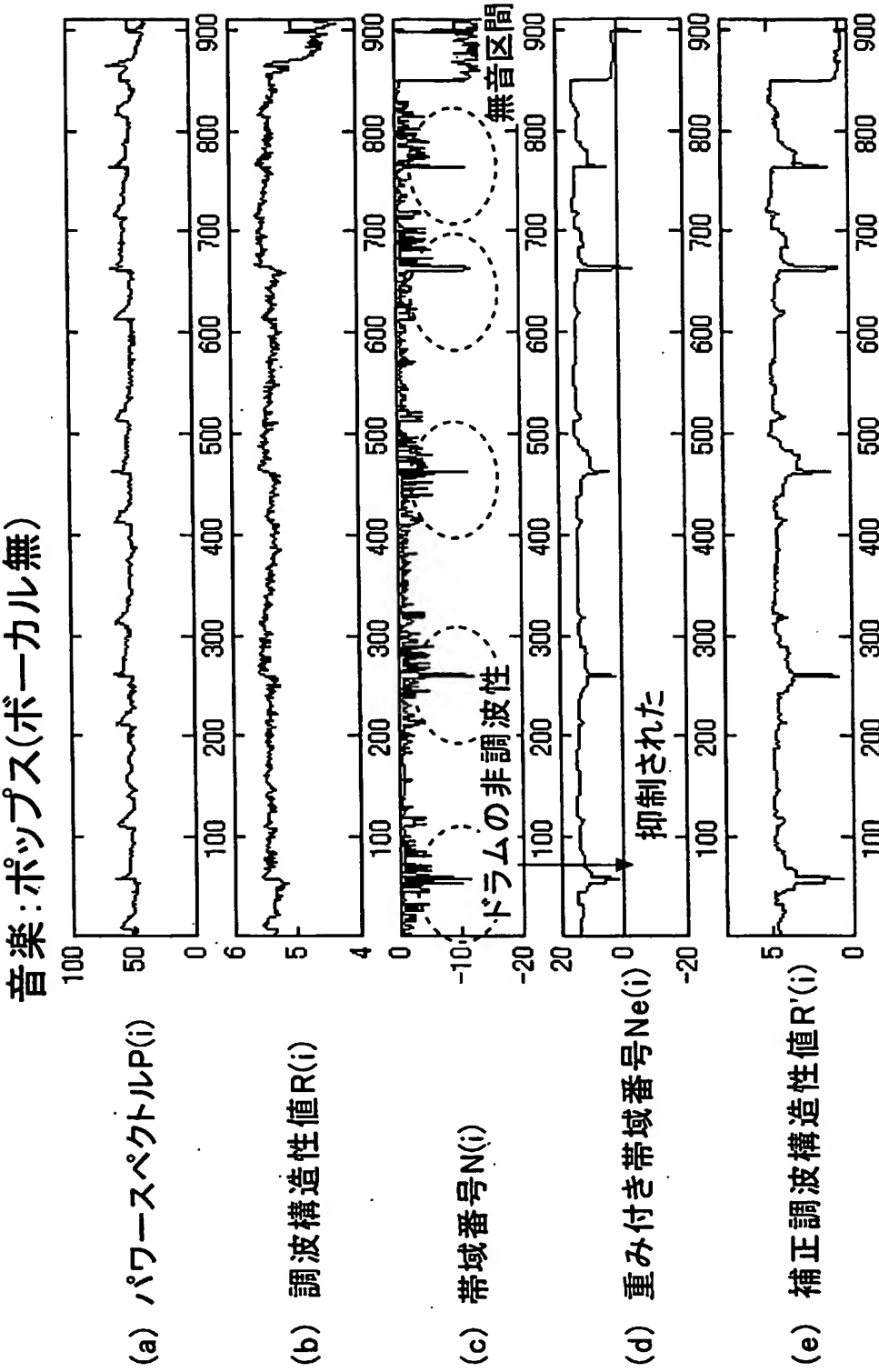


図16

音声：掃除機ノイズ00dB下

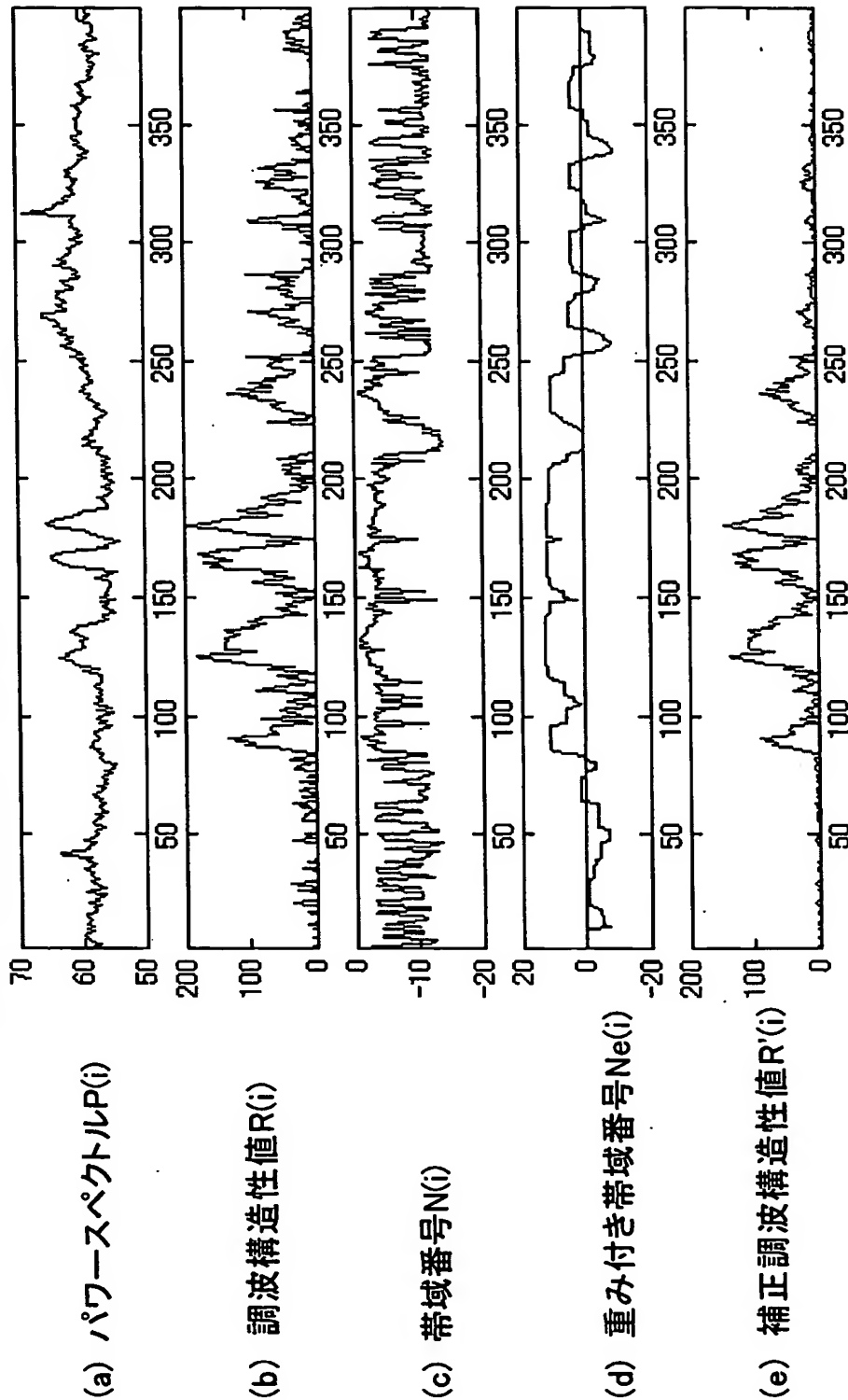


図17

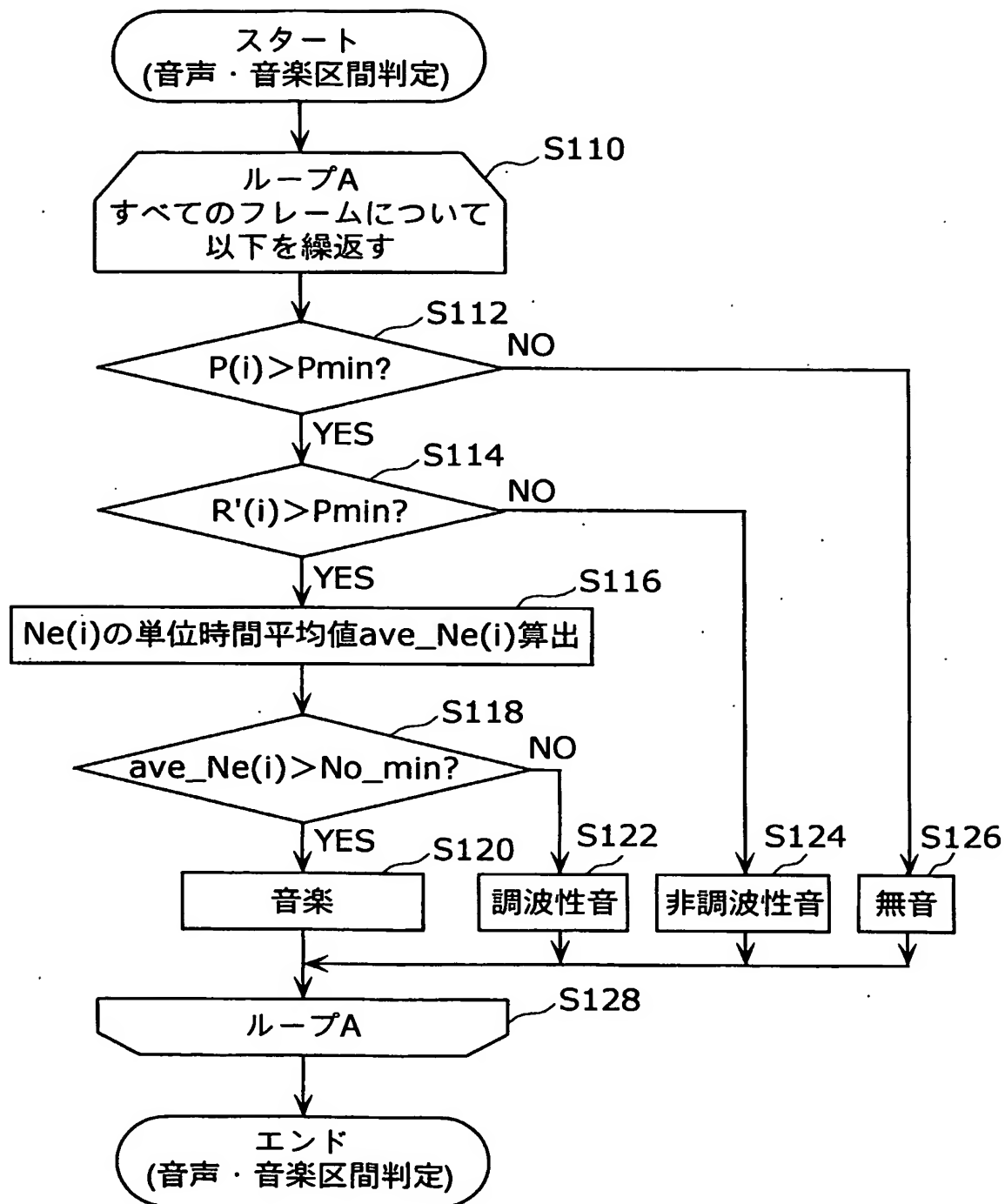


図18

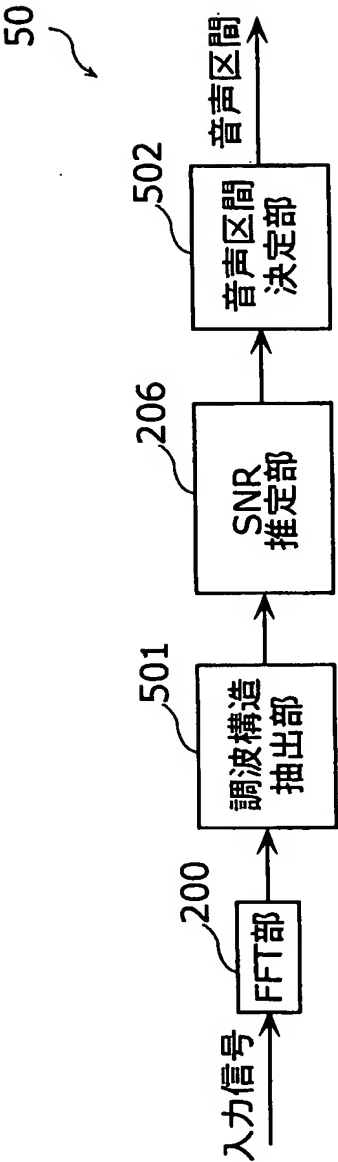


図19

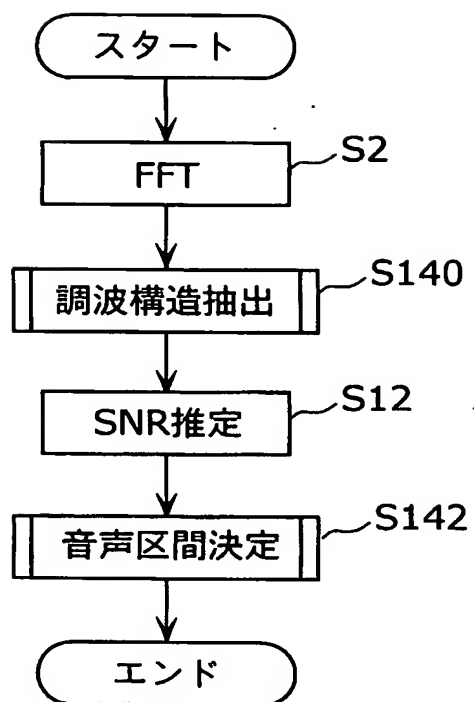


図20

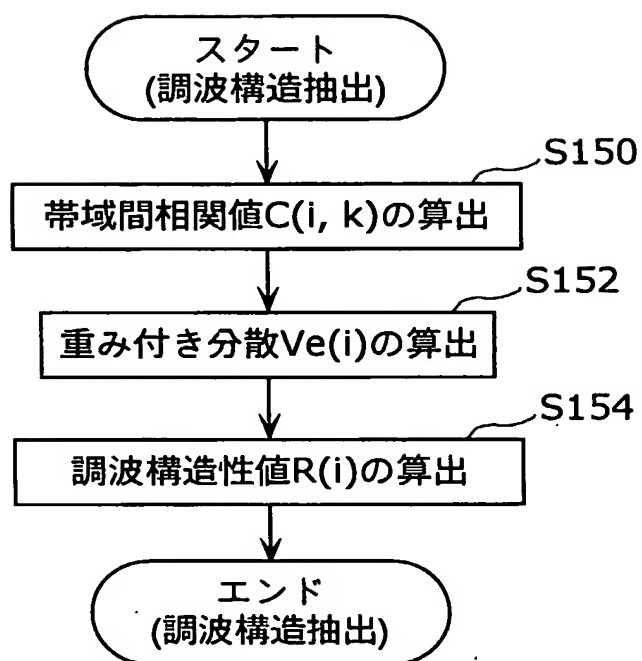


図21

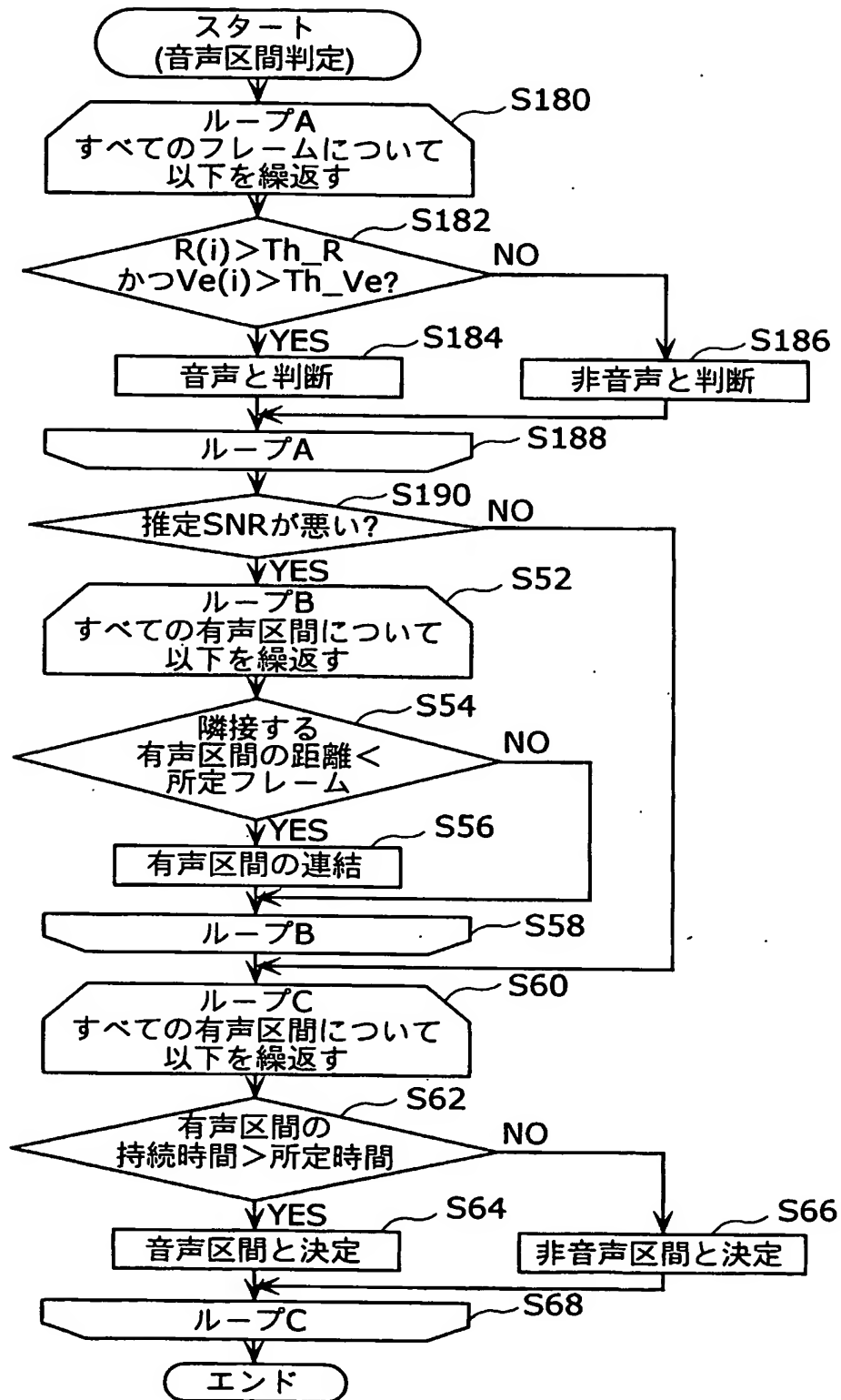


図22

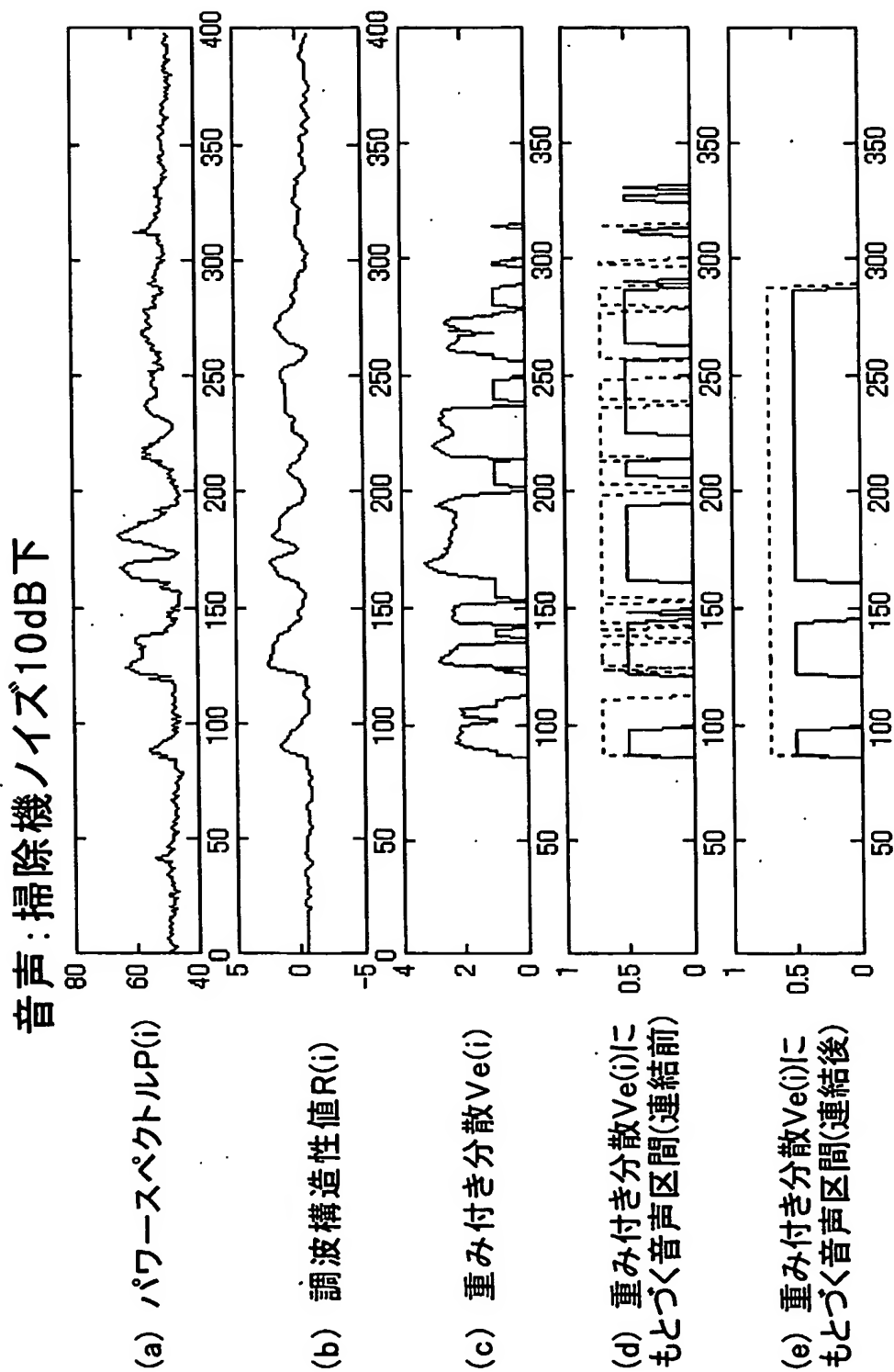


図23

音声：掃除機ノイズ40dB下

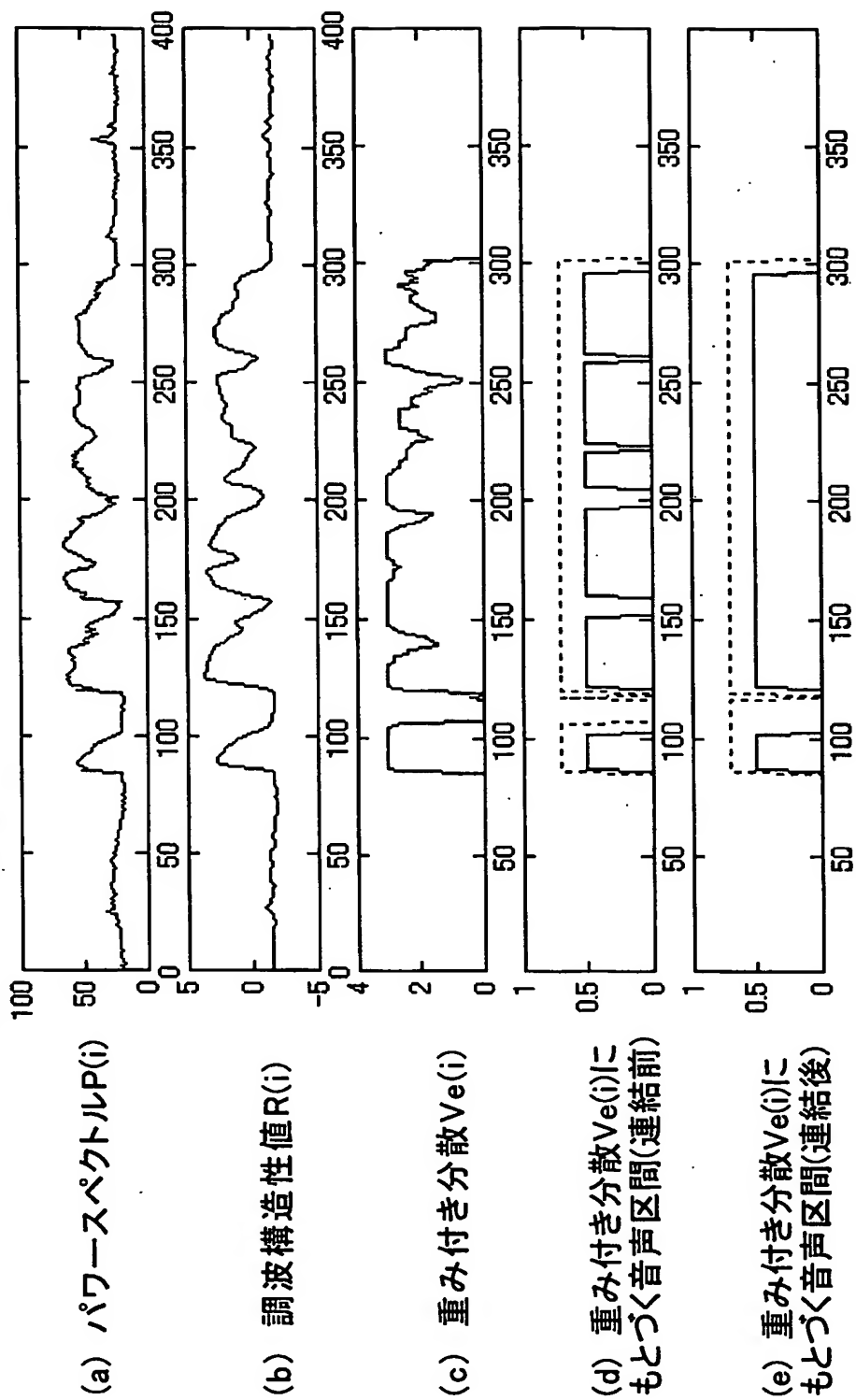


図24

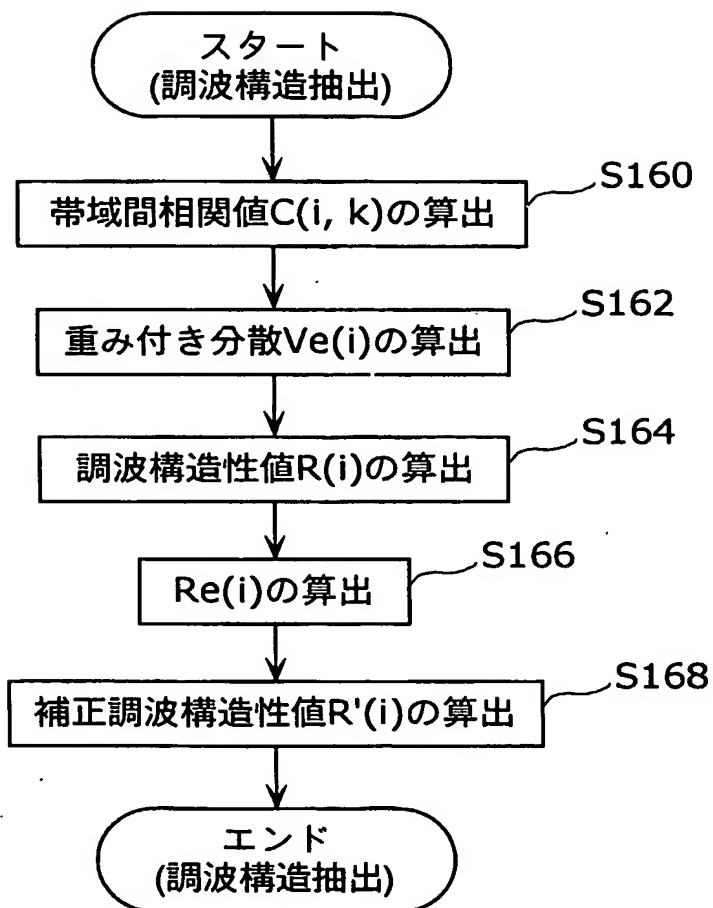


図25

音声: 掃除機ノイズ40dB下(ノイズ無)

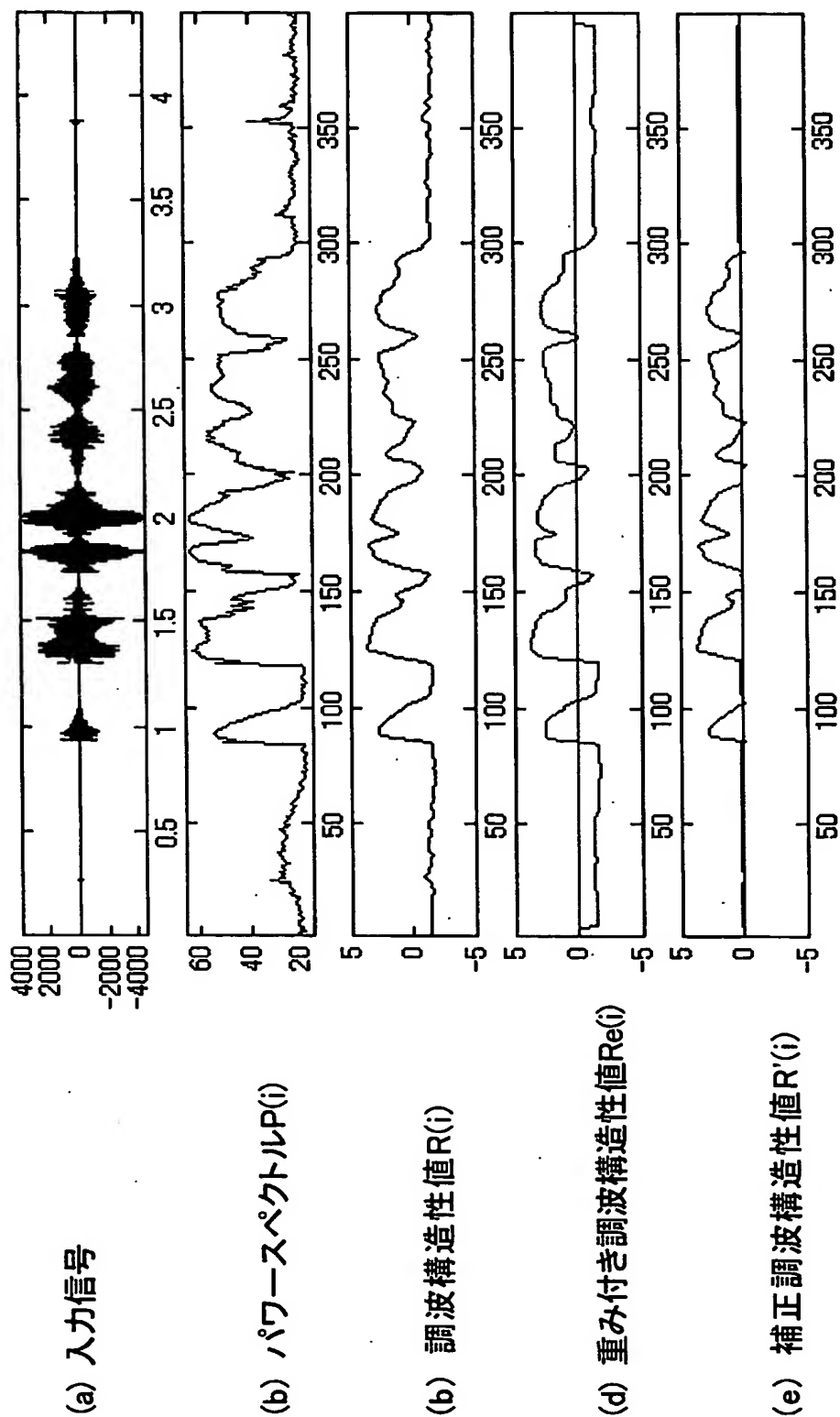


図26

音声: 掃除機ノイズ10dB下(ノイズ有)

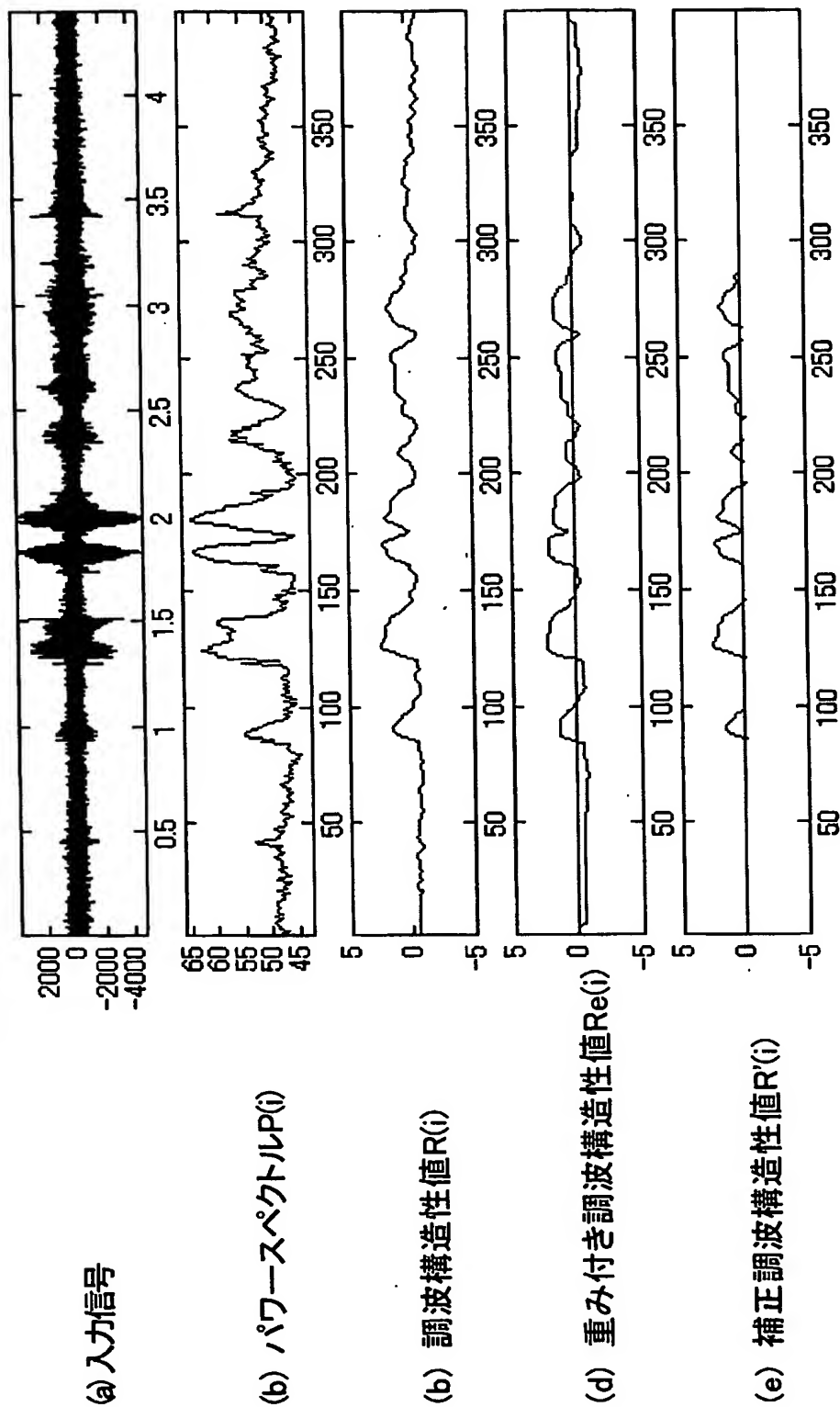


図27

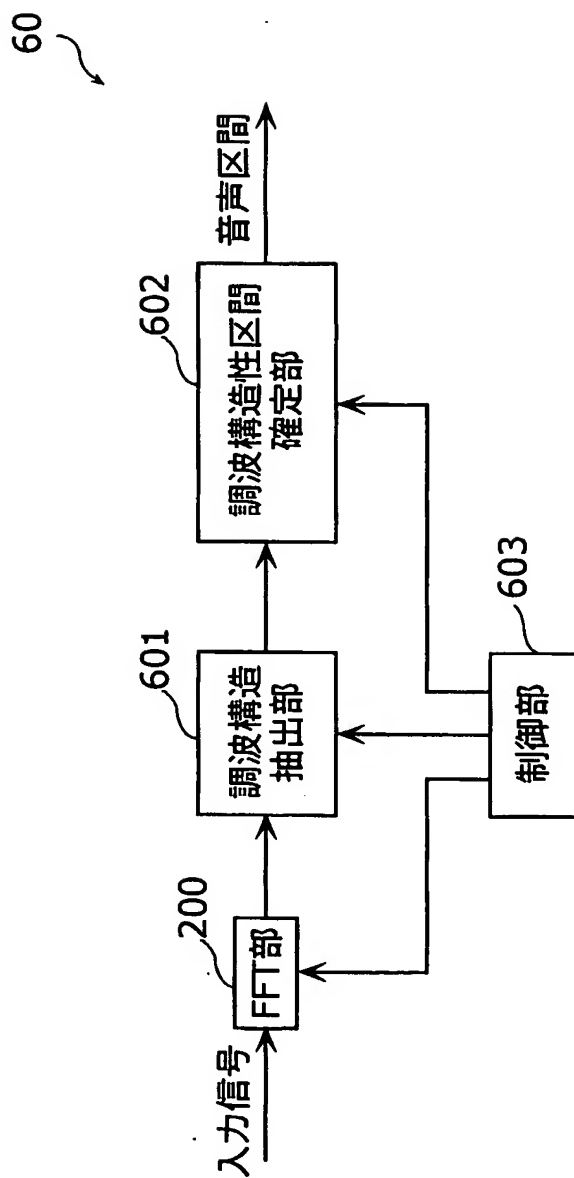


図28

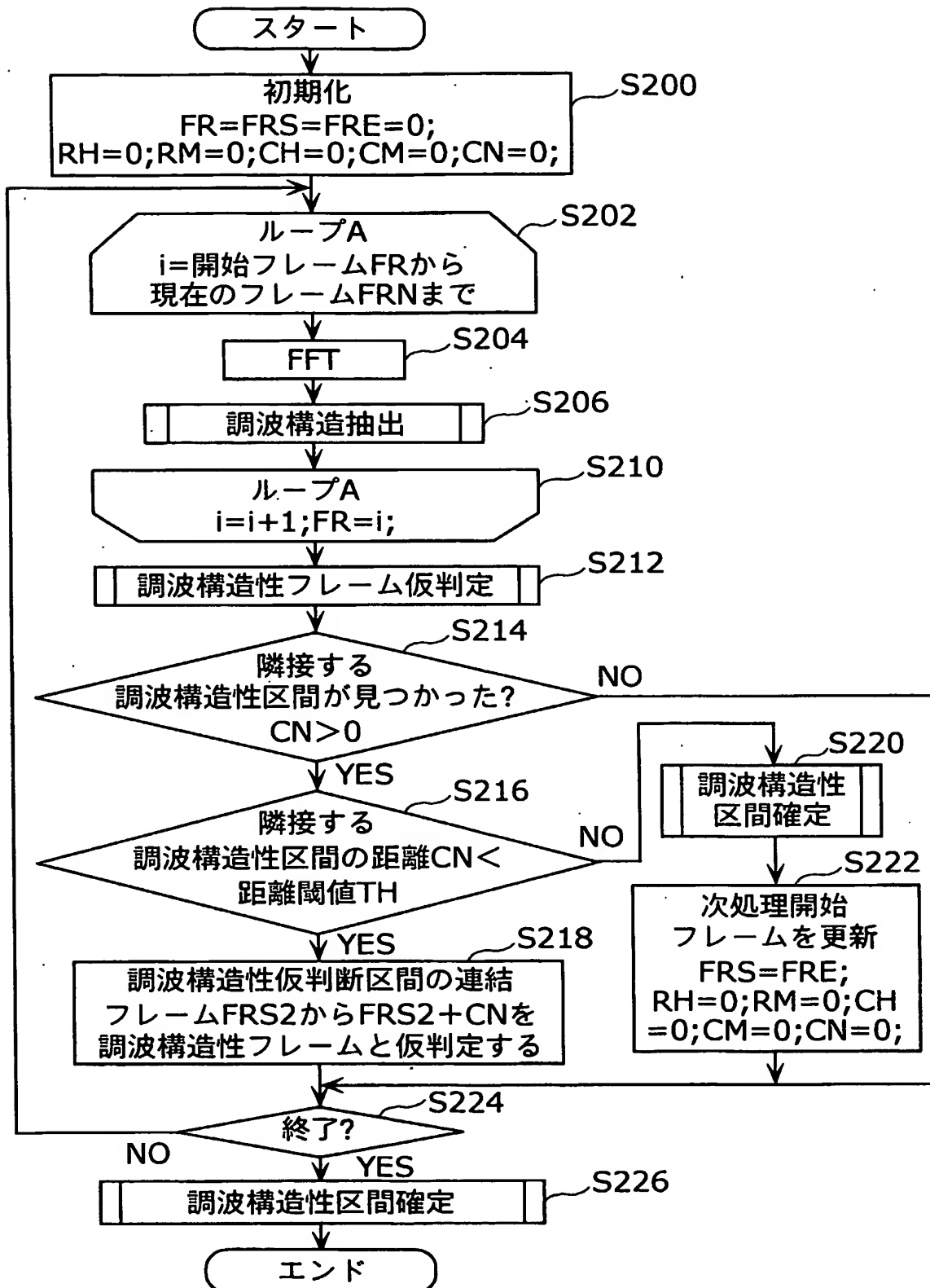


図29

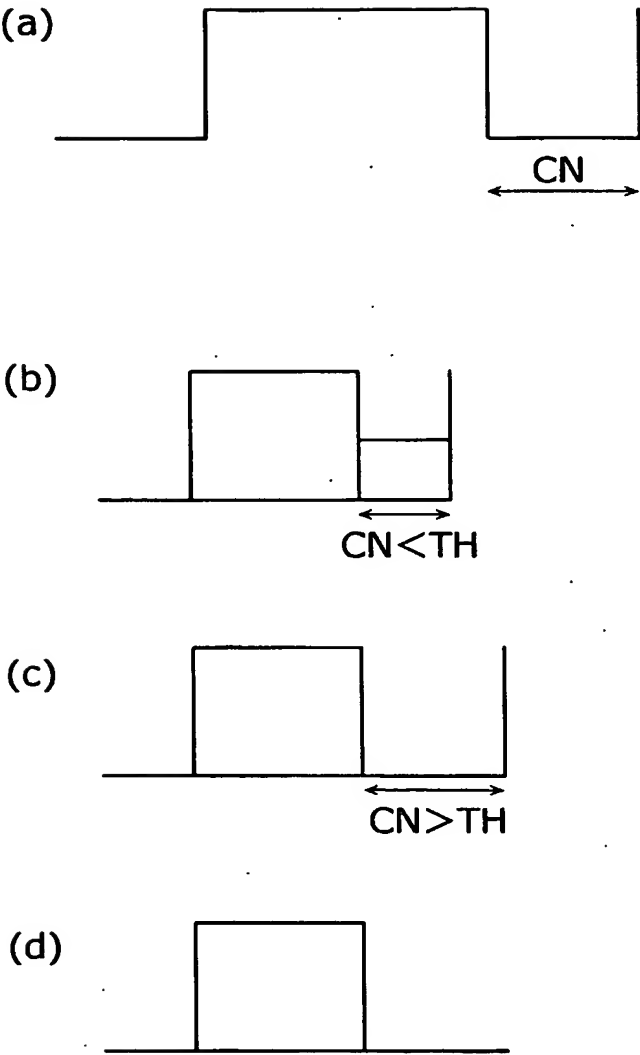


図30

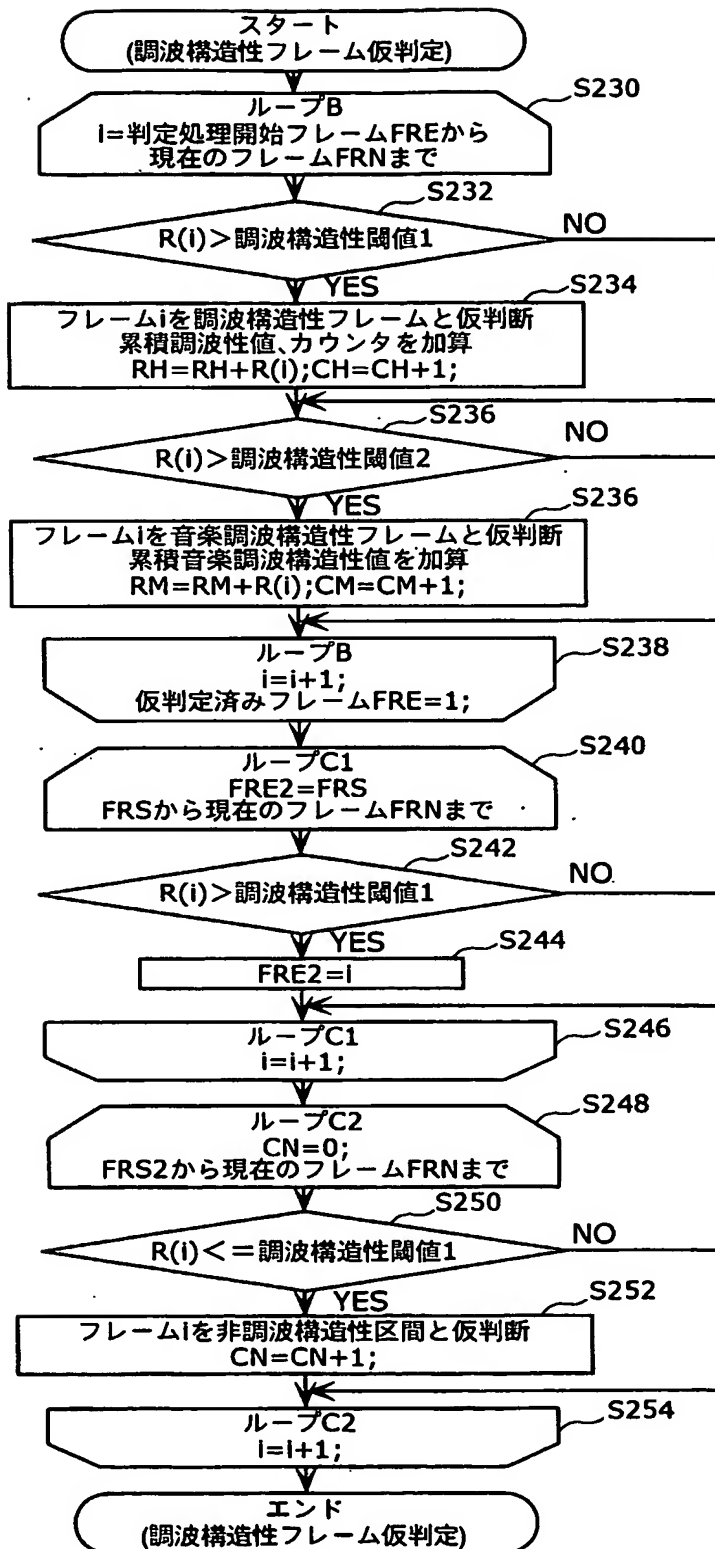


図31

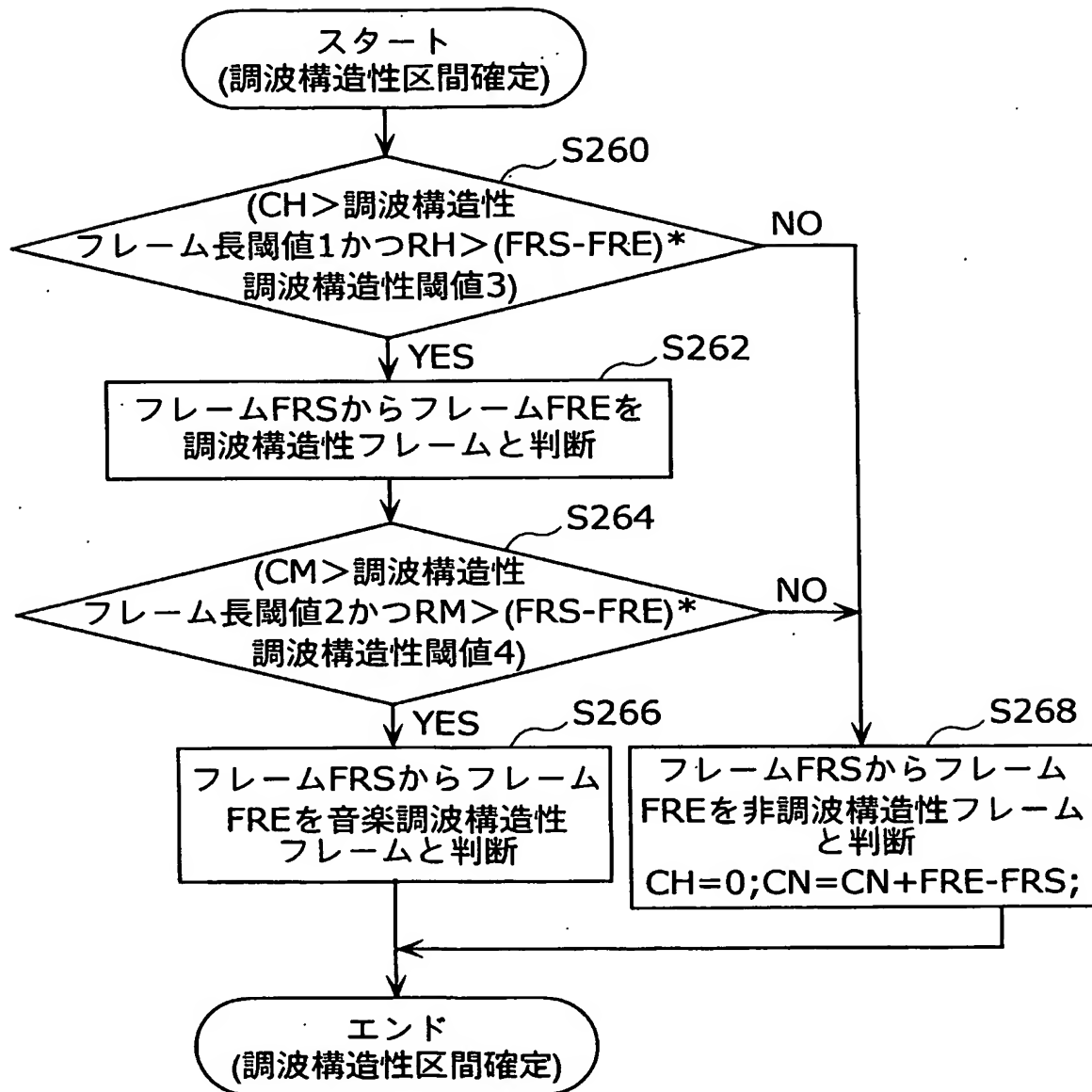
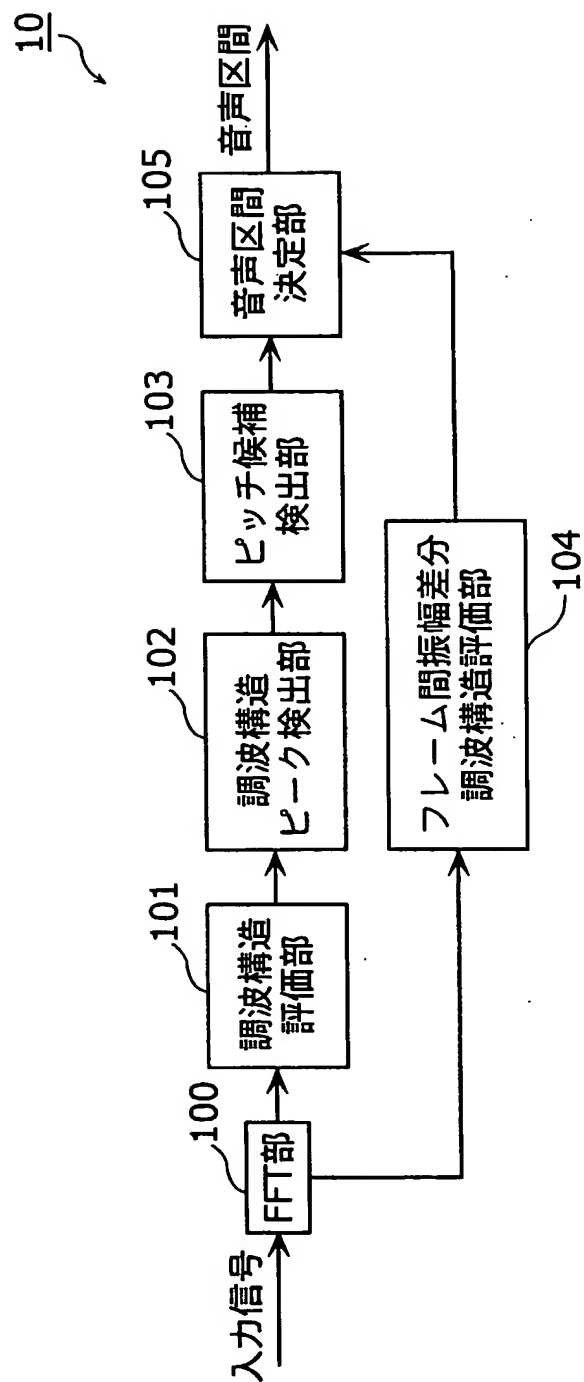


図32



INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2004/008051

A. CLASSIFICATION OF SUBJECT MATTER
Int.Cl⁷ G10L11/02, 06

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)
Int.Cl⁷ G10L11/02, 06Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched
Jitsuyo Shinan Koho 1922-1996 Toroku Jitsuyo Shinan Koho 1996-2004
Kokai Jitsuyo Shinan Koho 1971-2004 Jitsuyo Shinan Toroku Koho 1994-2004Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)
JSTPlus FILE (JOIS)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	JP 2002-162982 A (Matsushita Electric Industrial Co., Ltd.), 07 June, 2002 (07.06.02), Full text; Figs. 1 to 4 (Family: none)	1-17, 19-21
A	JP 2001-236085 A (Matsushita Electric Industrial Co., Ltd.), 31 August, 2001 (31.08.01), Full text; Figs. 1 to 15 (Family: none)	1-17, 19-21
A	JP 2001-222289 A (Yamaha Corp.), 17 August, 2001 (17.08.01), Full text; Figs. 1 to 7 (Family: none)	1-17, 19-21

☒ Further documents are listed in the continuation of Box C.☐ See patent family annex.

* Special categories of cited documents:

- "A" document defining the general state of the art which is not considered to be of particular relevance
- "E" earlier application or patent but published on or after the international filing date
- "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)
- "O" document referring to an oral disclosure, use, exhibition or other means
- "P" document published prior to the international filing date but later than the priority date claimed

- "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
- "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
- "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
- "&" document member of the same patent family

Date of the actual completion of the international search
17 September, 2004 (17.09.04)Date of mailing of the international search report
05 October, 2004 (05.10.04)Name and mailing address of the ISA/
Japanese Patent Office

Authorized officer

Facsimile No.

Telephone No.

INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2004/008051

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	JP 9-153769 A (Nippon Telegraph And Telephone Corp.), 10 June, 1997 (10.06.97), Full text; Figs. 1 to 6 (Family: none)	1-17,19-21
A	JP 11-143460 A (Nippon Telegraph And Telephone Corp.), 28 May, 1999 (28.05.99), Full text; Figs. 1 to 8 (Family: none)	1-17,19-21

INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2004/008051

Box No. II Observations where certain claims were found unsearchable (Continuation of item 2 of first sheet)

This international search report has not been established in respect of certain claims under Article 17(2)(a) for the following reasons:

1. ☐ Claims Nos.:
because they relate to subject matter not required to be searched by this Authority, namely:

2. ☒ Claims Nos.: 18
because they relate to parts of the international application that do not comply with the prescribed requirements to such an extent that no meaningful international search can be carried out, specifically:
Claim 18 describes "a harmonic-structured acoustic signal interval detection device as claimed in claim 28". However, no claim 28 exists in the scope of claims. Accordingly, claim 18 lacks the requirement of clarity within the meaning PCT Article 6.
3. ☐ Claims Nos.:
because they are dependent claims and are not drafted in accordance with the second and third sentences of Rule 6.4(a).

Box No. III Observations where unity of invention is lacking (Continuation of item 3 of first sheet)

This International Searching Authority found multiple inventions in this international application, as follows:

1. ☐ As all required additional search fees were timely paid by the applicant, this international search report covers all searchable claims.
2. ☐ As all searchable claims could be searched without effort justifying an additional fee, this Authority did not invite payment of any additional fee.
3. ☐ As only some of the required additional search fees were timely paid by the applicant, this international search report covers only those claims for which fees were paid, specifically claims Nos.:

4. ☐ No required additional search fees were timely paid by the applicant. Consequently, this international search report is restricted to the invention first mentioned in the claims; it is covered by claims Nos.:

Remark on Protest

- ☐ The additional search fees were accompanied by the applicant's protest.
☐ No protest accompanied the payment of additional search fees.

A. 発明の属する分野の分類 (国際特許分類 (IPC))

Int. Cl⁷ G10L11/02, 06

B. 調査を行った分野

調査を行った最小限資料 (国際特許分類 (IPC))

Int. Cl⁷ G10L11/02, 06

最小限資料以外の資料で調査を行った分野に含まれるもの

日本国実用新案公報	1922-1996年
日本国公開実用新案公報	1971-2004年
日本国登録実用新案公報	1996-2004年
日本国実用新案登録公報	1994-2004年

国際調査で使用した電子データベース (データベースの名称、調査に使用した用語)

JSTPlusファイル (JOIS)

C. 関連すると認められる文献

引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求の範囲の番号
A	J P 2002-162982 A (松下電器産業株式会社) 2002. 06. 07、全文、第1-4図 (ファミリーなし)	1-17, 19-21
A	J P 2001-236085 A (松下電器産業株式会社) 2001. 08. 31、全文、第1-15図 (ファミリーなし)	1-17, 19-21
A	J P 2001-222289 A (ヤマハ株式会社) 2001. 08. 17、全文、第1-7図 (ファミリーなし)	1-17, 19-21

☒ C欄の続きにも文献が列挙されている。☐ パテントファミリーに関する別紙を参照。

* 引用文献のカテゴリー

「A」 特に関連のある文献ではなく、一般的技術水準を示すもの
「E」 国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの
「L」 優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す)
「O」 口頭による開示、使用、展示等に言及する文献
「P」 国際出願日前で、かつ優先権の主張の基礎となる出願

の日の後に公表された文献

「T」 国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの
「X」 特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの
「Y」 特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの
「&」 同一パテントファミリー文献

国際調査を完了した日

17. 09. 2004

国際調査報告の発送日

05.10.2004

国際調査機関の名称及びあて先

日本国特許庁 (ISA/J P)
郵便番号100-8915
東京都千代田区霞が関三丁目4番3号

特許庁審査官 (権限のある職員)

南 義明

5C

9381

電話番号 03-3581-1101 内線 3541

C (続き). 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求の範囲の番号
A	JP 9-153769 A (日本電信電話株式会社) 1997. 06. 10、全文、第1-6図 (ファミリーなし)	1-17, 19-21
A	JP 11-143460 A (日本電信電話株式会社) 1999. 05. 28、全文、第1-8図 (ファミリーなし)	1-17, 19-21

第Ⅱ欄 請求の範囲の一部の調査ができないときの意見 (第1ページの2の続き)

法第8条第3項 (PCT 17条(2)(a)) の規定により、この国際調査報告は次の理由により請求の範囲の一部について作成しなかった。

1. ☐ 請求の範囲 _____ は、この国際調査機関が調査をすることを要しない対象に係るものである。
つまり、
2. ☒ 請求の範囲 18 は、有意義な国際調査をすることができる程度まで所定の要件を満たしていない国際出願の部分に係るものである。つまり、
「請求の範囲 28 項に記載の調波構造的音響信号区間検出装置」と記載する請求項 18 は、請求項 28 項が請求の範囲に存在しない故、PCT 第6条における明確性の要件を欠いている。
3. ☐ 請求の範囲 _____ は、従属請求の範囲であって PCT 規則 6.4(a) の第2文及び第3文の規定に従って記載されていない。

第Ⅲ欄 発明の単一性が欠如しているときの意見 (第1ページの3の続き)

次に述べるようにこの国際出願に二以上の発明があるところの国際調査機関は認めた。

1. ☐ 出願人が必要な追加調査手数料をすべて期間内に納付したので、この国際調査報告は、すべての調査可能な請求の範囲について作成した。
2. ☐ 追加調査手数料を要求するまでもなく、すべての調査可能な請求の範囲について調査することができたので、追加調査手数料の納付を求めなかった。
3. ☐ 出願人が必要な追加調査手数料を一部のみしか期間内に納付しなかったため、この国際調査報告は、手数料の納付のあった次の請求の範囲のみについて作成した。
4. ☐ 出願人が必要な追加調査手数料を期間内に納付しなかったため、この国際調査報告は、請求の範囲の最初に記載されている発明に係る次の請求の範囲について作成した。

追加調査手数料の異議の申立てに関する注意

- ☐ 追加調査手数料の納付と共に出願人から異議申立てがあった。
☐ 追加調査手数料の納付と共に出願人から異議申立てがなかった。